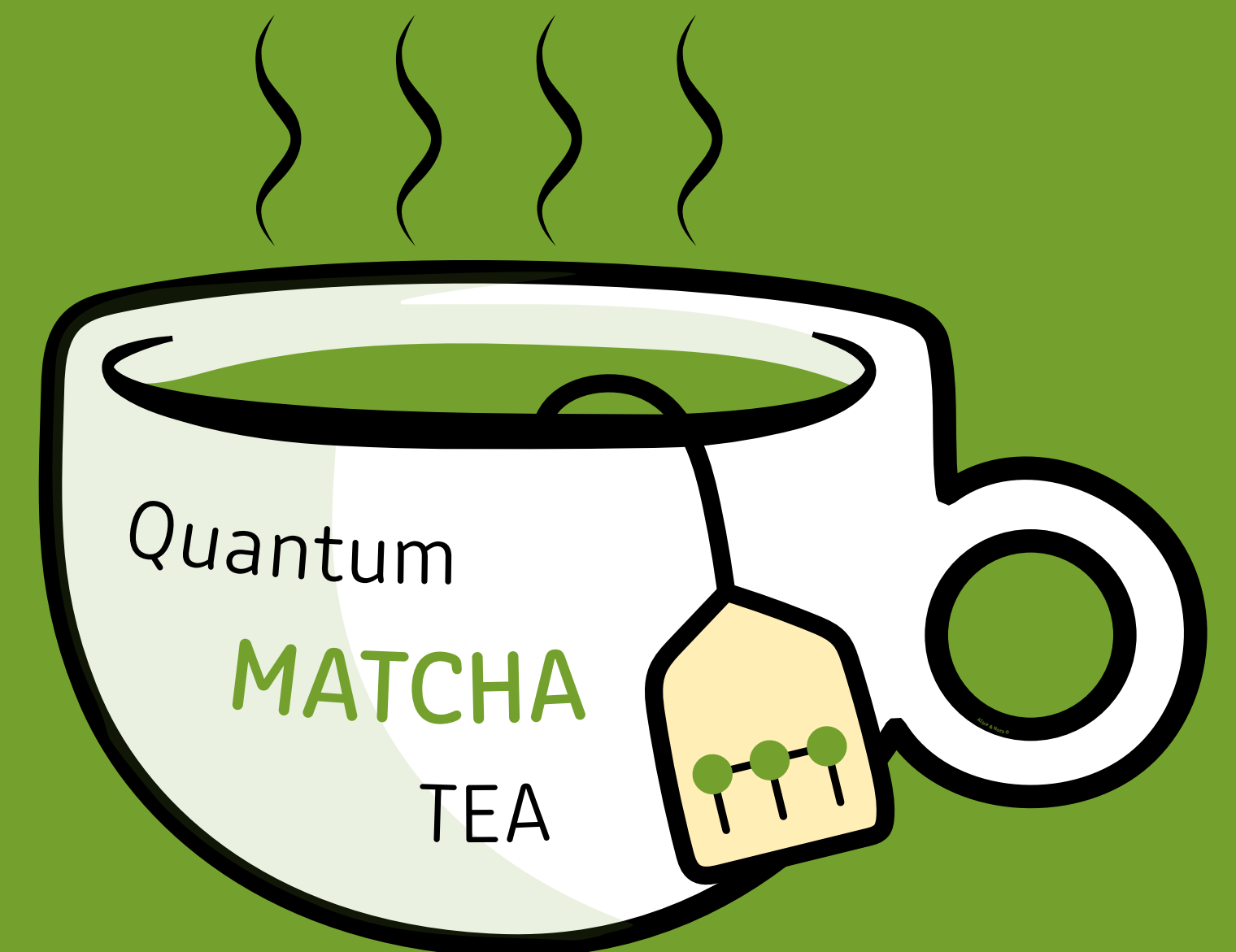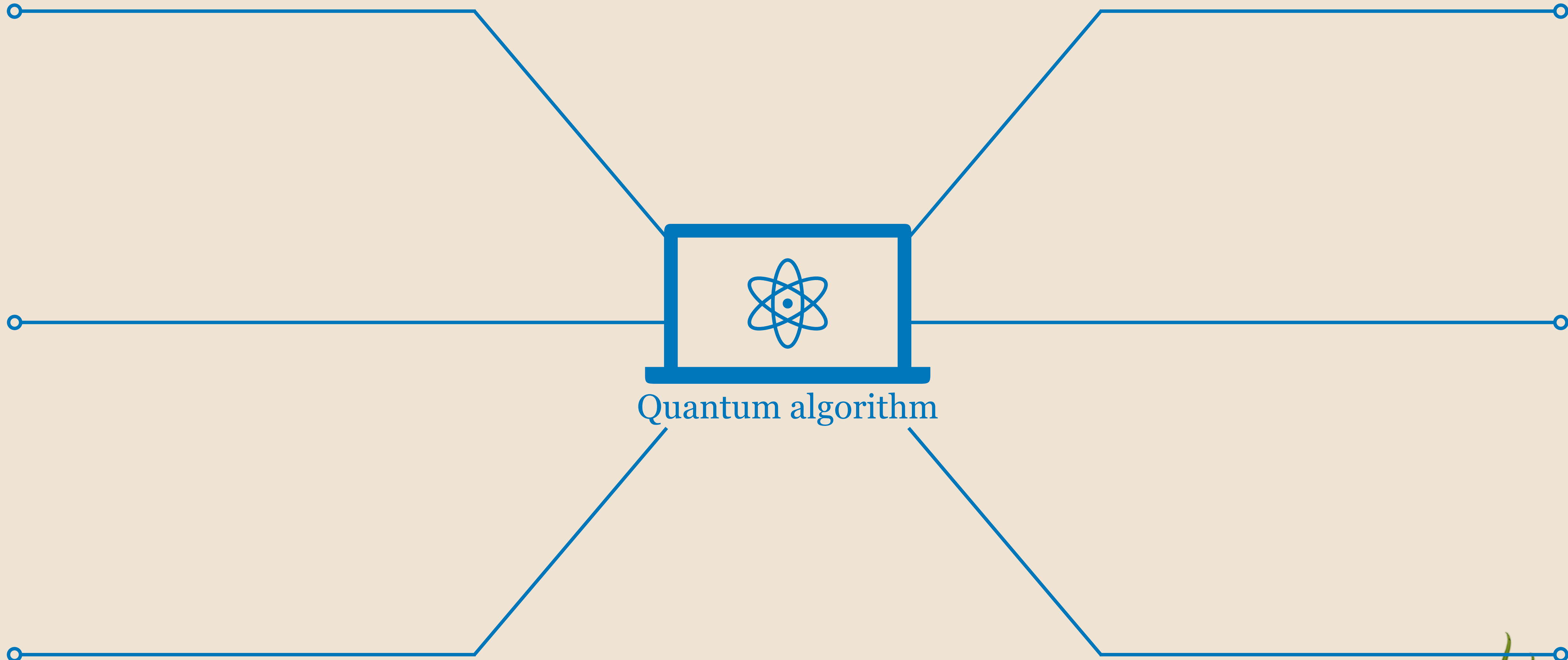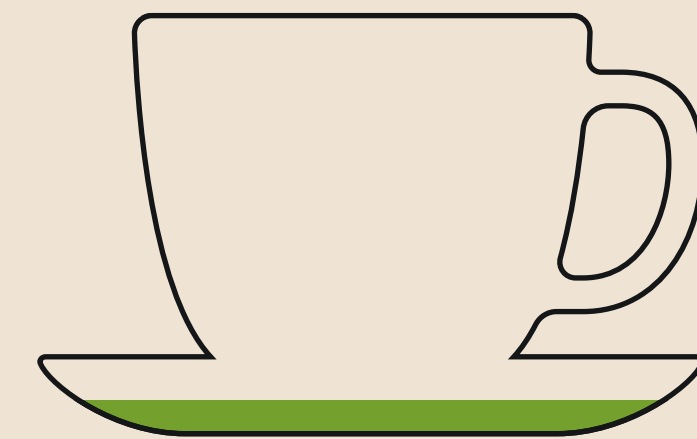# Quantum Matcha Tea

*An efficient matrix product state simulator for quantum circuits*

**Marco Ballarin**
**Università degli studi di Padova**

Quantum
MATCHA
TEA

# Running quantum algorithms

Quantum algorithm

# Running quantum algorithms

+ Real hardware
- Noisy
- Limited number of qubits
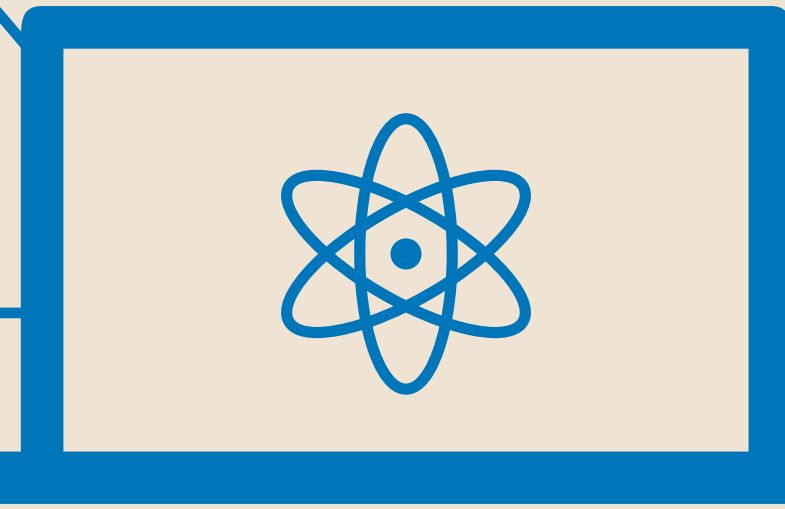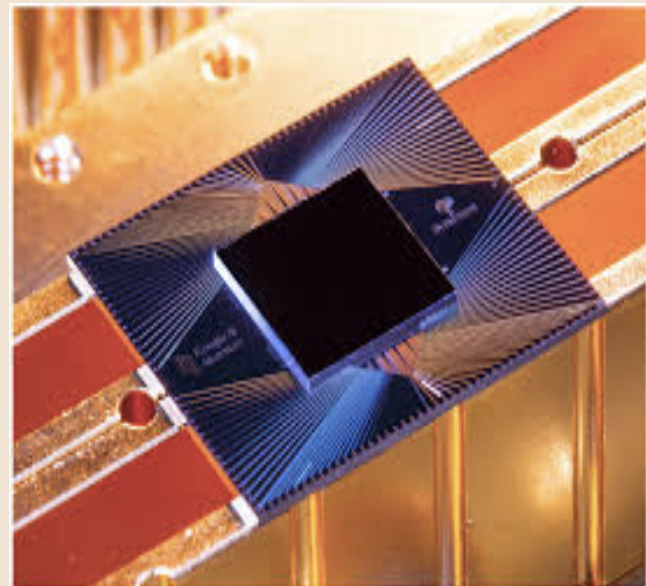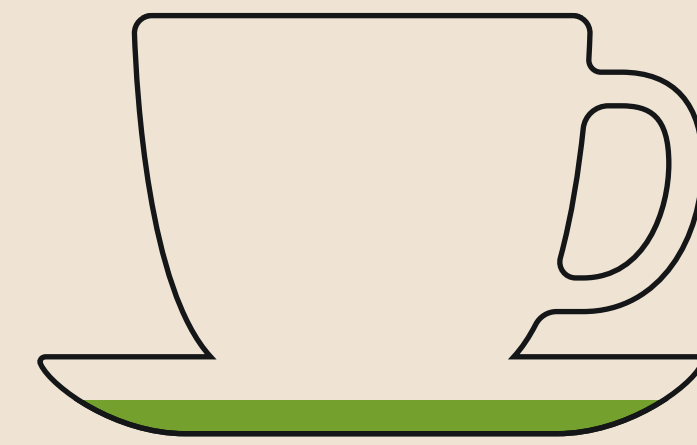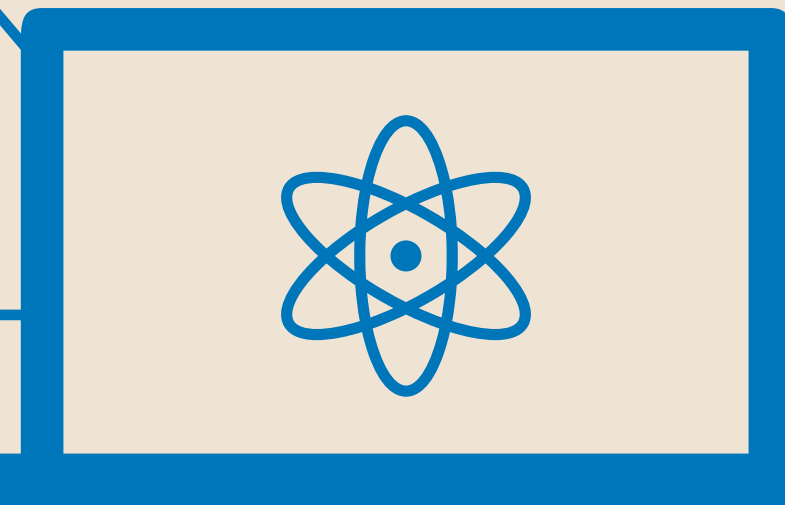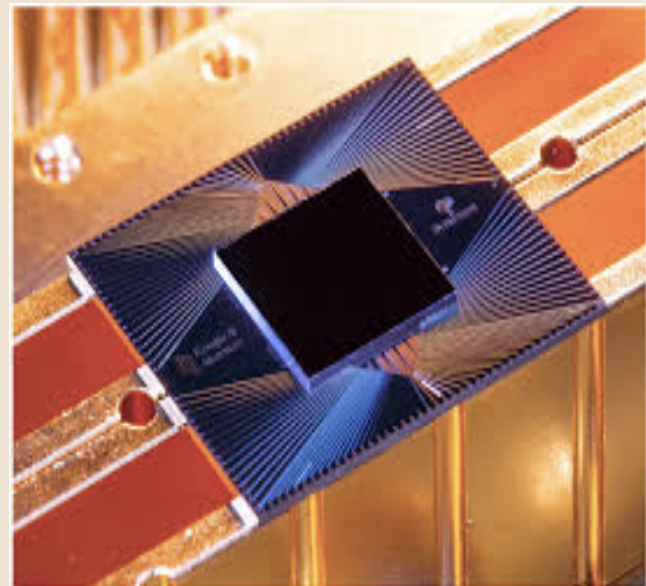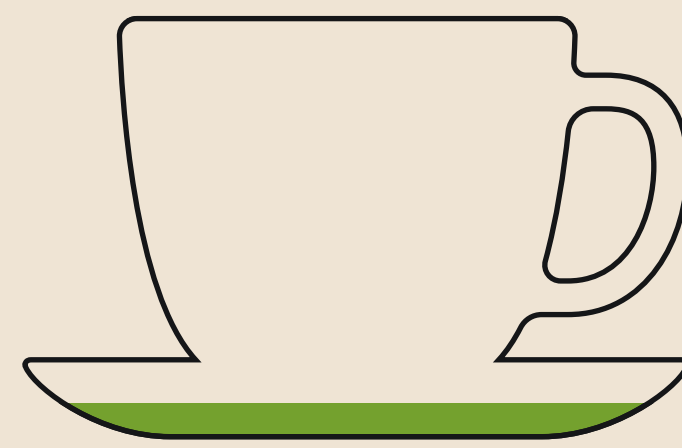
Quantum hardware

Quantum algorithm

# Running quantum algorithms



+ Real hardware
- Noisy
- Limited number of qubits

Quantum hardware

High # of qubits +
Flexibility (observables) -
Depth of the circuit -

cuQuantum

Quantum algorithm

# Running quantum algorithms



Quantum hardware
+ Real hardware
- Noisy
- Limited number of qubits

cuQuantum
High # of qubits +
Flexibility (observables) -
Depth of the circuit -

Quantum algorithm

Exact simulator
+ Access to exact state
- Limited number of qubits

2

# Running quantum algorithms



Quantum hardware

+ Real hardware
- Noisy
- Limited number of qubits

cuQuantum

High # of qubits +
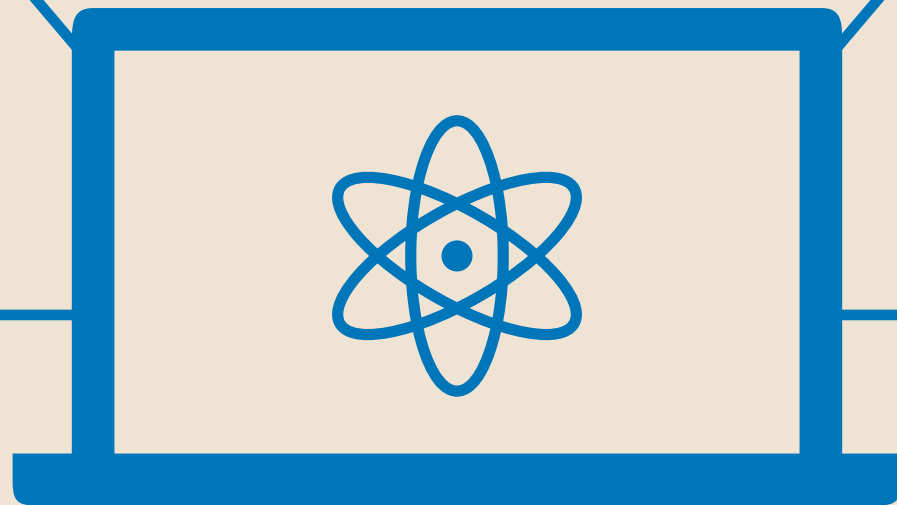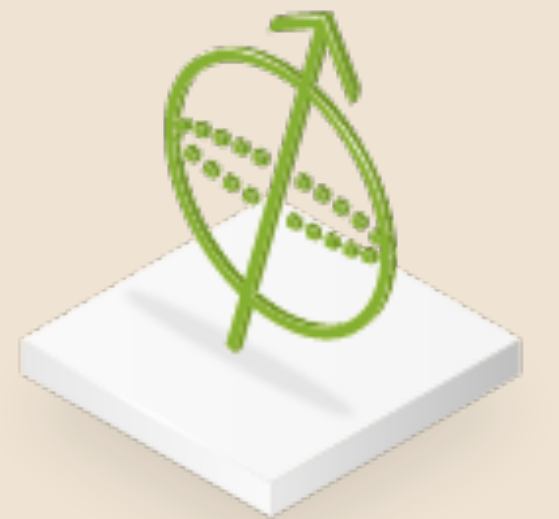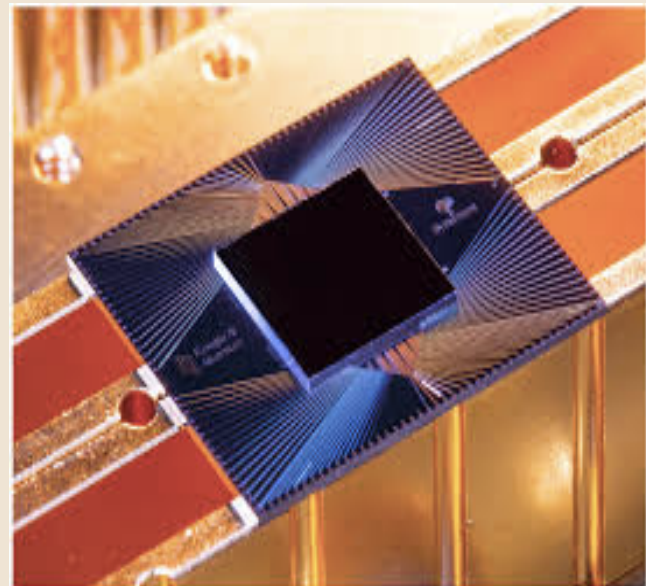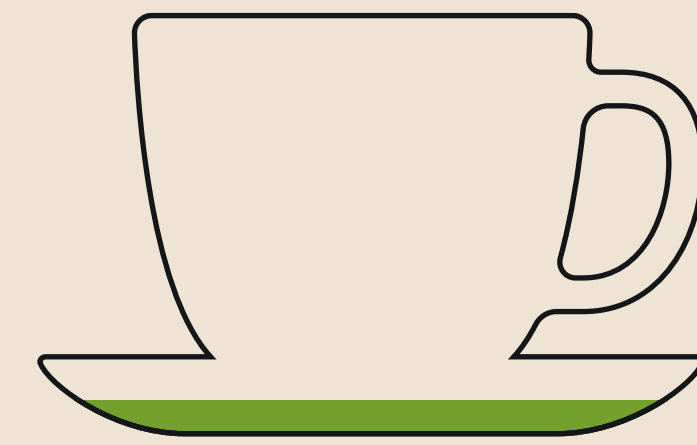Flexibility (observables) -
Depth of the circuit -

Quantum algorithm

Exact simulator
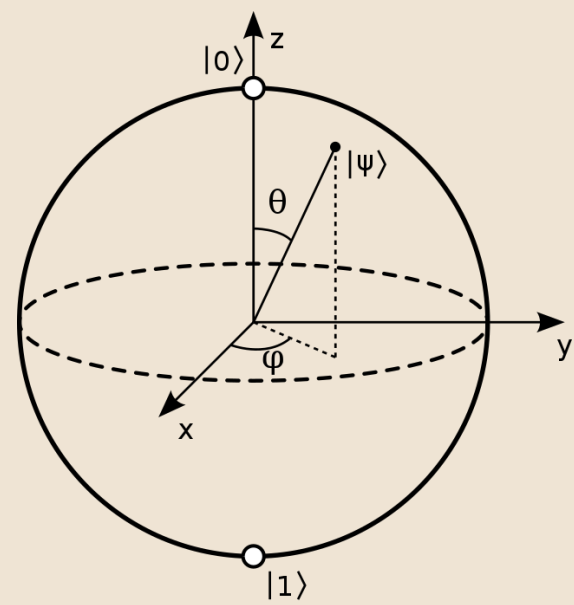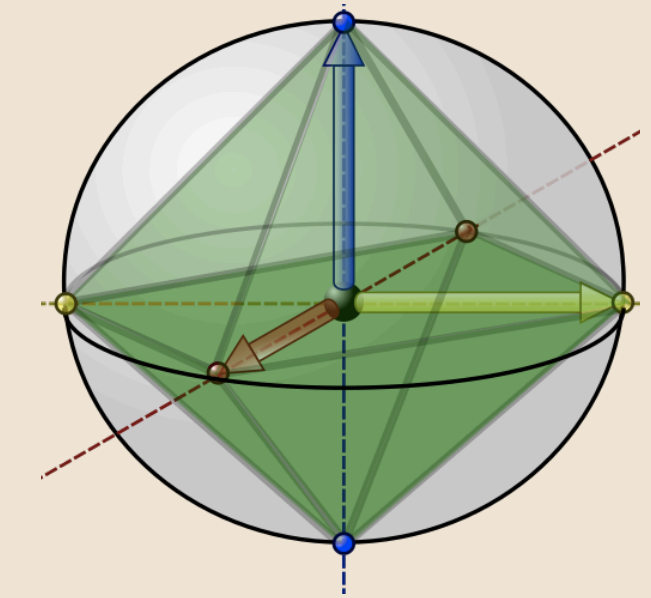
+ Access to exact state
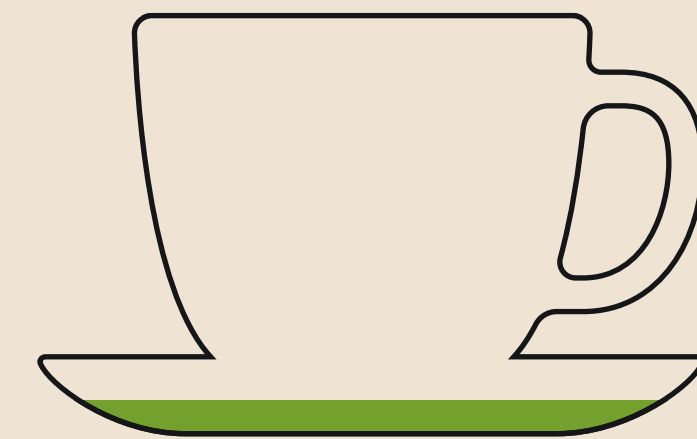- Limited number of qubits

Clifford simulator

High # of qubits +
Flexibility (# of T gates) -

# Running quantum algorithms



**Quantum hardware**
+ Real hardware
- Noisy
- Limited number of qubits

**cuQuantum**
High # of qubits +
Flexibility (observables) -
Depth of the circuit -

**Quantum algorithm**

**Exact simulator**
+ Access to exact state
- Limited number of qubits

**Tensor Network simulator**
+ High # of qubits
- Flexibility (entanglement)
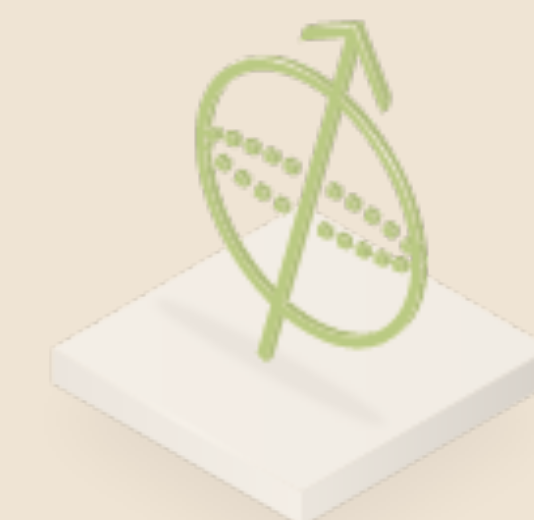
**Clifford simulator**
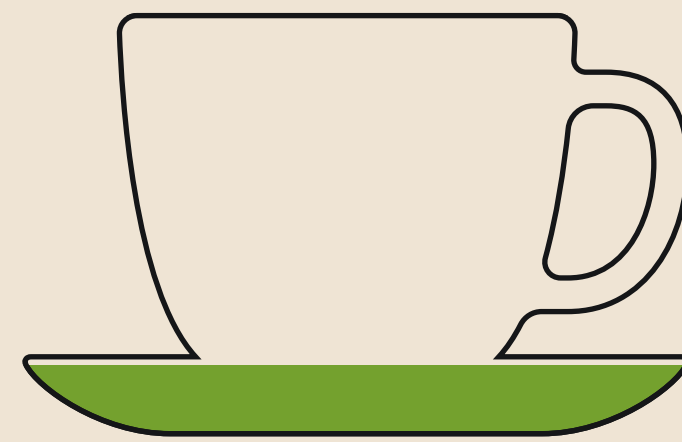High # of qubits +
Flexibility (# of T gates) -

# Why tensor networks

$$dim(\mathcal{H}) = 2^n$$

We can represent a
subset efficiently

?

# Why tensor networks

$dim(\mathcal{H}) = 2^n$

?

We can represent a
subset efficiently

# Why tensor networks

$dim(\mathcal{H}) = 2^n$

? 

We can represent a
subset efficiently

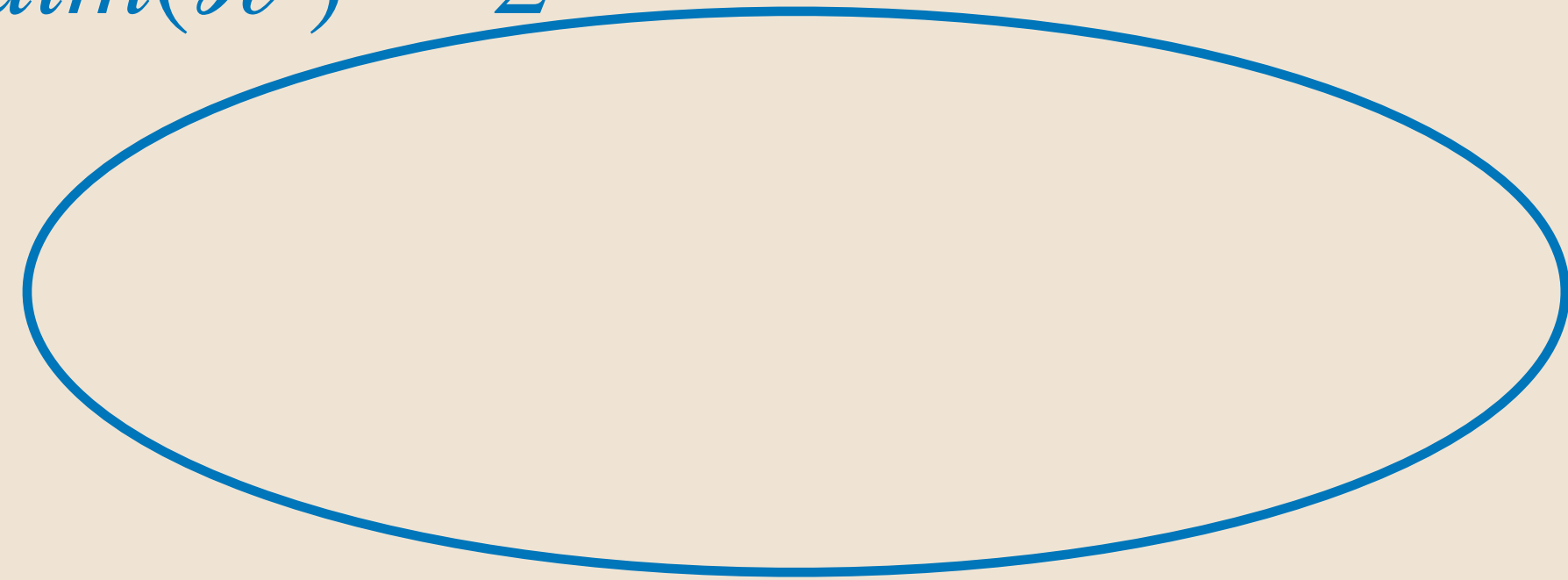$$|\psi\rangle = \sum_{\alpha=1}^{\chi} \boxed{|A_\alpha\rangle} \overset{\lambda_\alpha}{\rule{2cm}{0.4pt}} \boxed{|B_\alpha\rangle}$$

# Why tensor networks
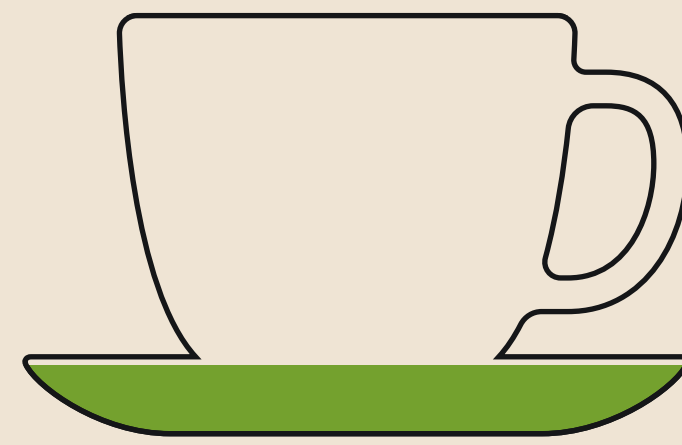
$dim(\mathcal{H}) = 2^n$

?

We can represent a subset efficiently

Tensor networks compress the quantum correlations between subsystems $\Rightarrow$ **compress entanglement**

$$|\psi\rangle = \sum_{\alpha=1}^{\chi}$$

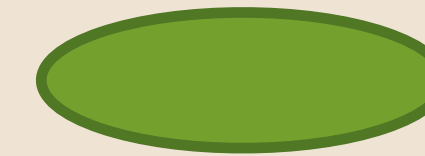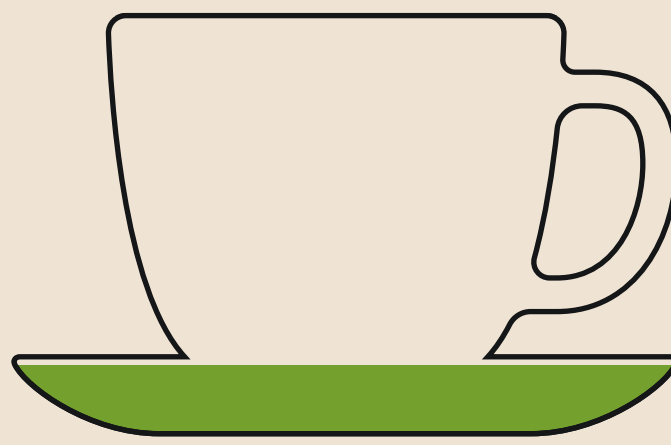$|A_\alpha\rangle$ — $\lambda_\alpha$ — $|B_\alpha\rangle$

# Why tensor networks

$dim(\mathcal{H}) = 2^n$

?

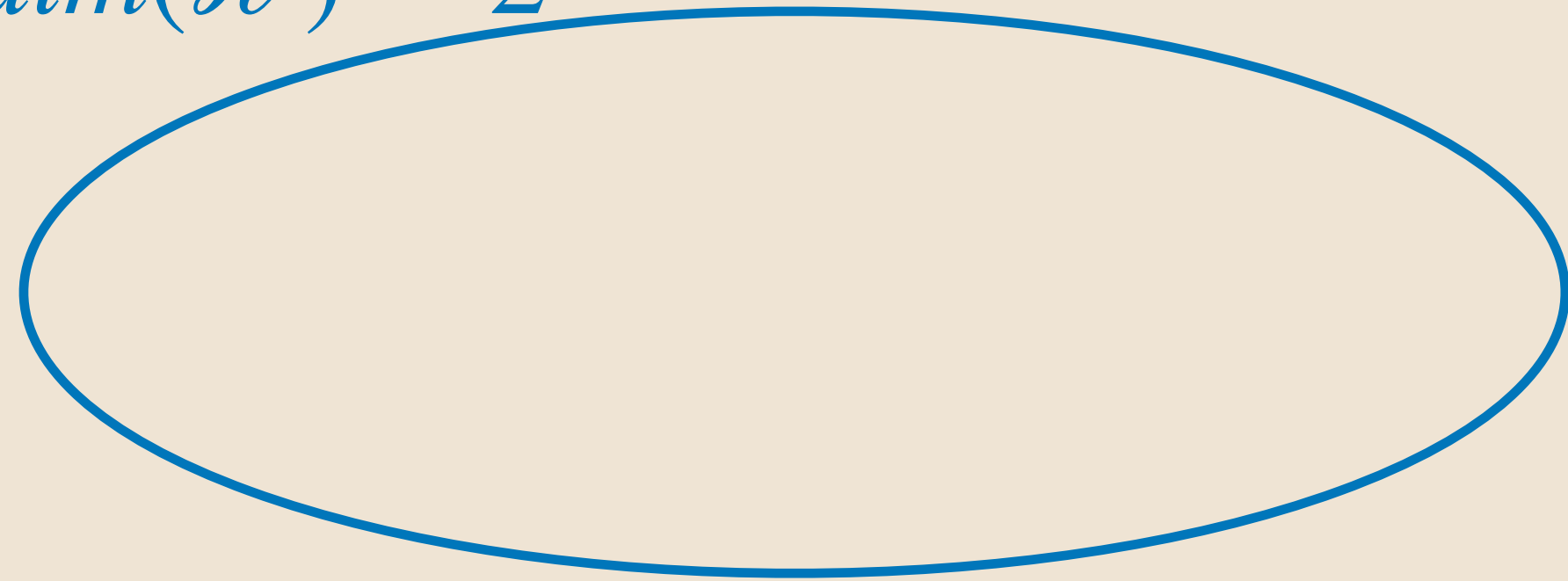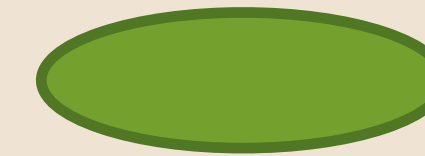We can represent a subset efficiently

Tensor networks compress the quantum correlations between subsystems ⇒ **compress entanglement**

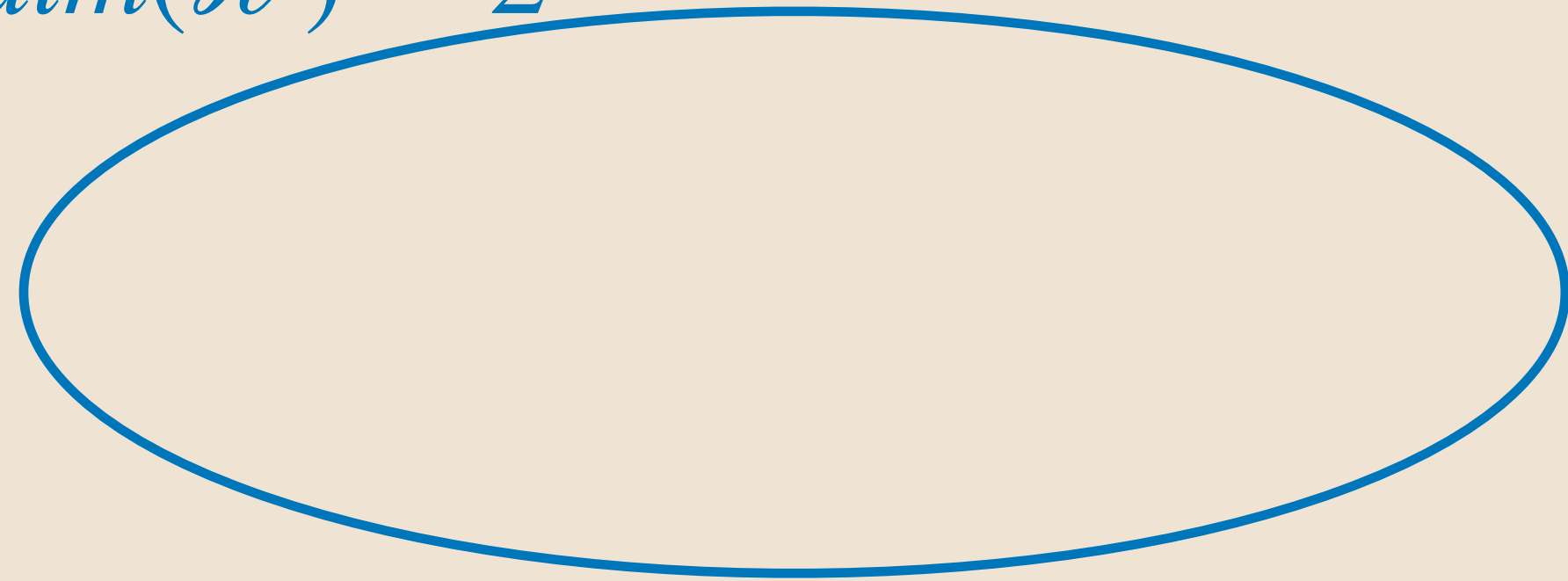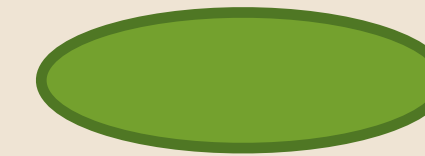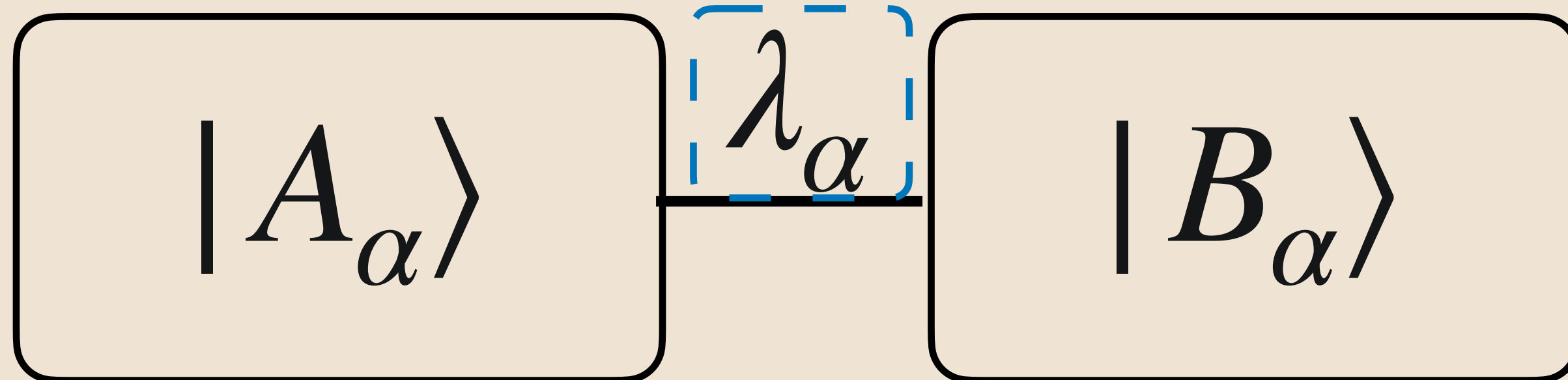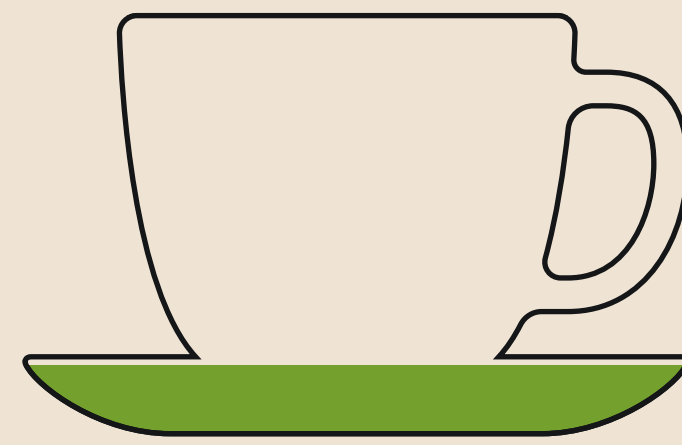$$|\psi\rangle = \sum_{\alpha=1}^{\chi} |A_\alpha\rangle \, \lambda_\alpha \, |B_\alpha\rangle$$

Only keep highest $\chi$ Schmidt values

# Matrix product states



$\chi$

Memory requirements

$$O(2^n) \rightarrow O(2n\chi^2)$$

# Matrix product states

Each tensor (ball) encodes the state of a qubit

$\chi$

Memory requirements

$$O(2^n) \rightarrow O(2n\chi^2)$$

# Matrix product states

Each tensor (ball) encodes
the state of a qubit

Bonds encode entanglement
between qubits

Memory requirements

$\chi$

$$O(2^n) \rightarrow O(2n\chi^2)$$

4

# Matrix product states

Each tensor (ball) encodes
the state of a qubit

Bonds encode entanglement
between qubits

Memory requirements

$$O(2^n) \rightarrow O(2n\chi^2)$$

$\chi$

State evolution through
quantum circuit



4

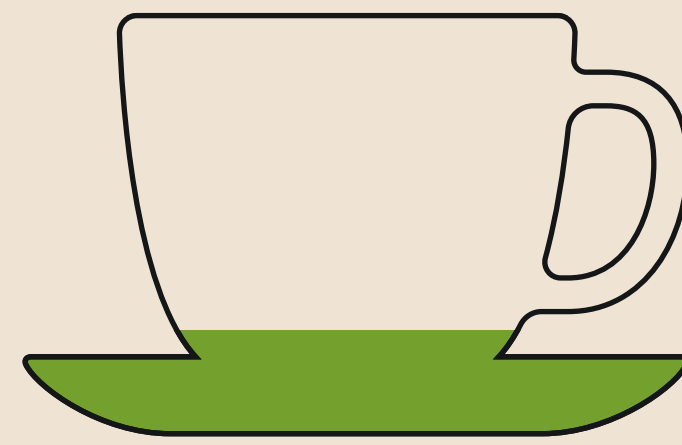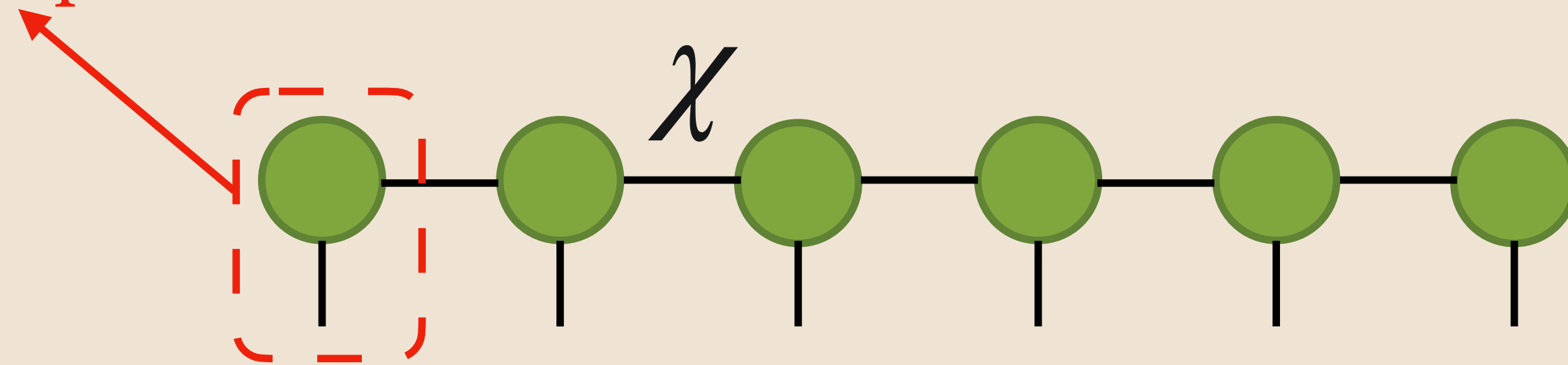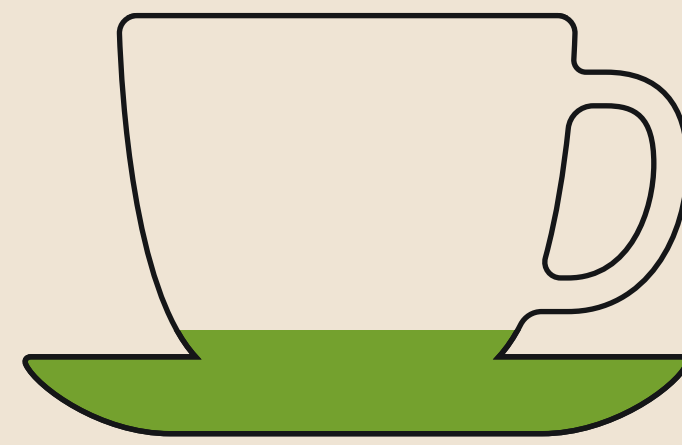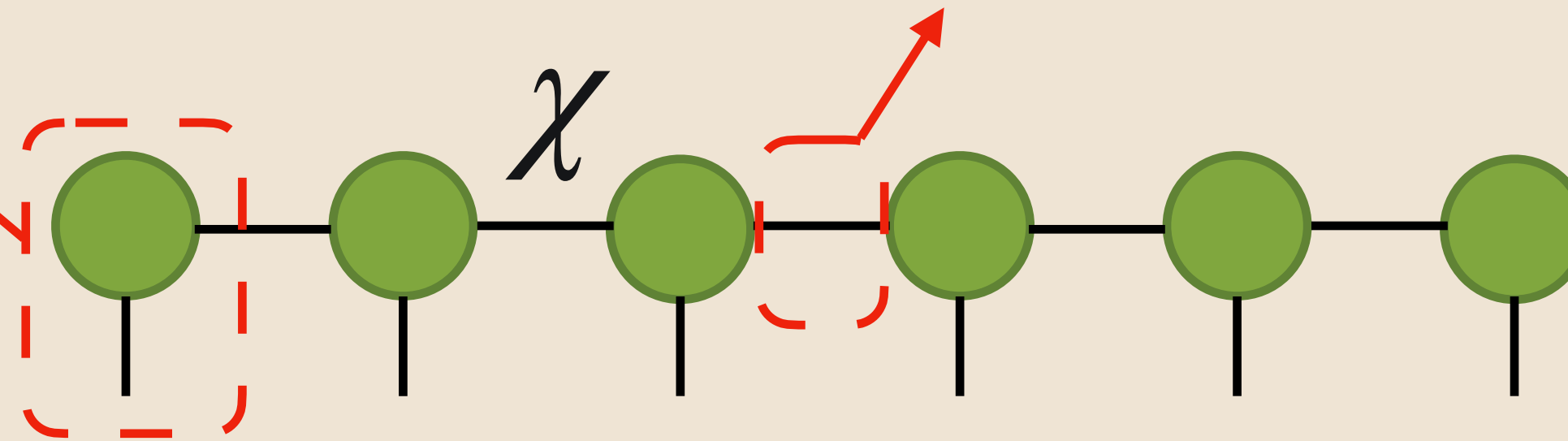# **Matrix product states**



Each tensor (ball) encodes
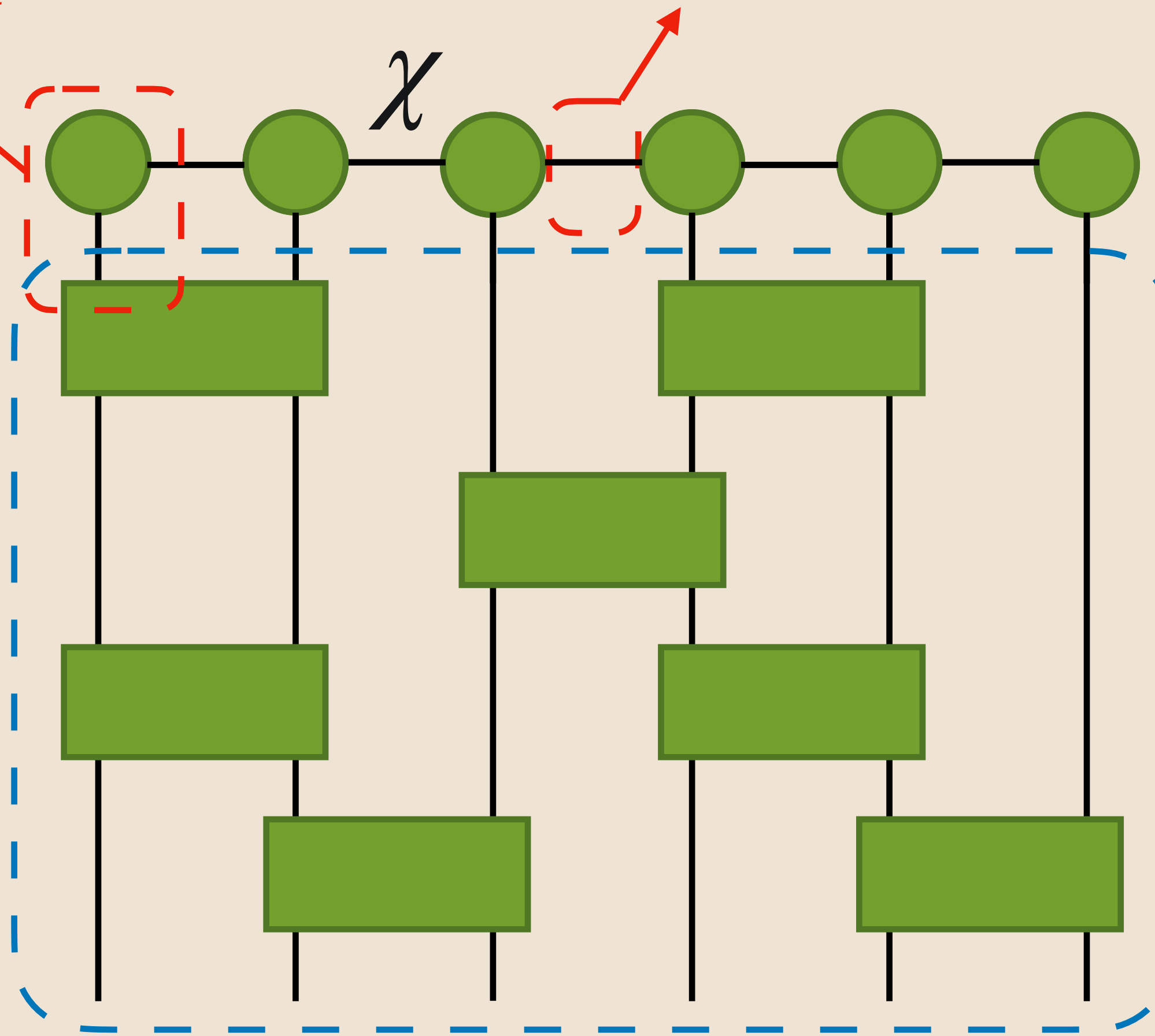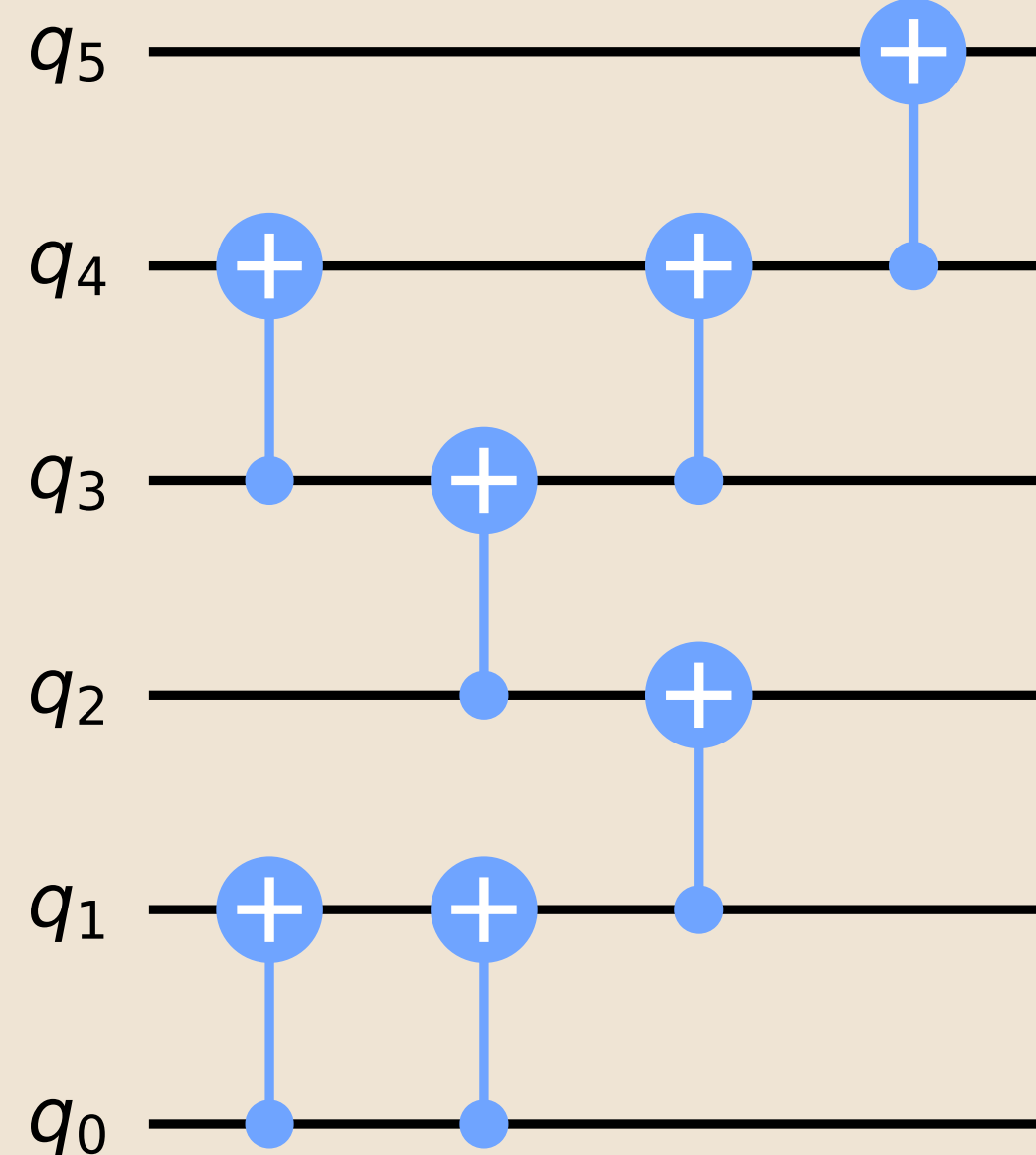the state of a qubit

Bonds encode entanglement
between qubits

$\chi$

Memory requirements

$$O(2^n) \to O(2n\chi^2)$$

State evolution through
quantum circuit

$q_5$

$q_4$

$q_3$

$q_2$

$q_1$

$q_0$

MPS SIMULATIONS ARE
NOT LIMITED BY THE
NUMBER OF QUBITS BUT
BY THE ENTANGLEMENT

4

# Quantum TEA distribution

**Tensor network** ← **T** **E** **A** → **Applications**

**Emulator**

# Quantum TEA distribution

Tensor network ← T | E | A → Applications

Emulator

Quantum tea leaves: **Utility**

Quantum matcha tea: **quantum circuit HPC simulations**

Quantum red tea: **tensor handling**

Quantum chai tea: **AI and ML with tensor networks**

Quantum green tea: **Schrödinger equation solution for many-body states**

# Quantum TEA distribution

Tensor network ← **T** **E** **A** → Applications

Emulator

**Public!**

Quantum tea leaves: **Utility**

Quantum matcha tea: **quantum circuit HPC simulations**

Quantum red tea: **tensor handling**

Quantum chai tea: **AI and ML with tensor networks**

Quantum green tea: **Schrödinger equation solution for many-body states**

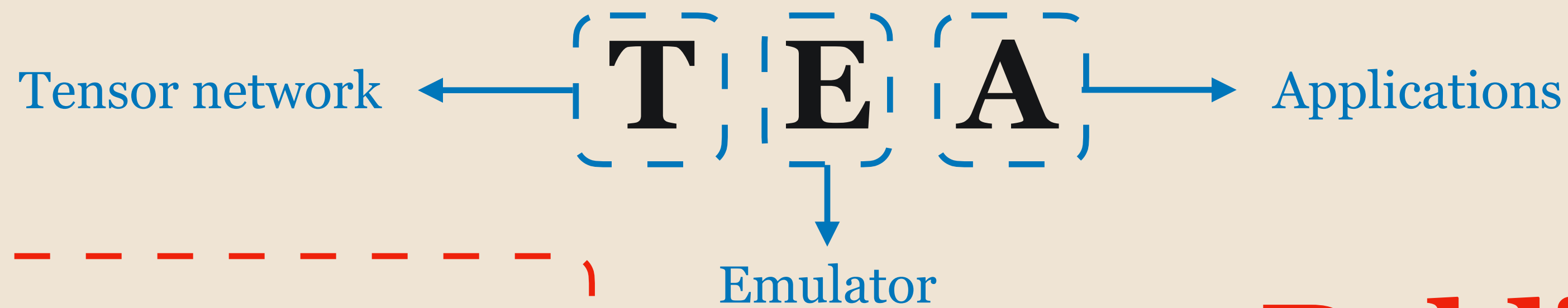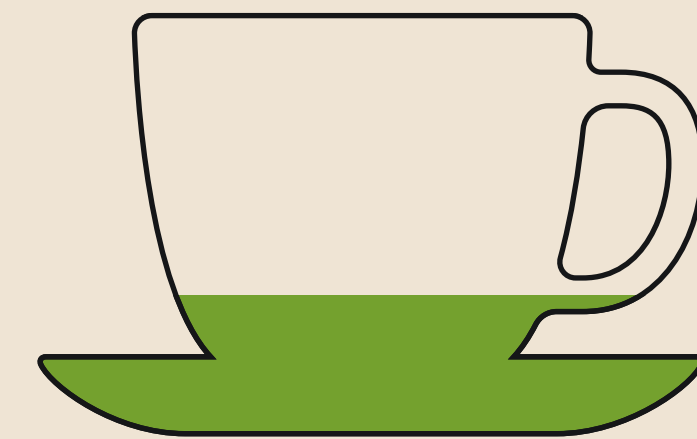# Quantum Matcha Tea workflow

Quantum circuit

Observables

Python interface, definition
of the problem

# Quantum Matcha Tea workflow

Quantum circuit

Observables

Matrix product state
simulator

# Quantum Matcha Tea workflow

Quantum circuit

Observables

Matrix product state simulator

NumPy

CuPy

F

Serial CPU
Multinode MPI CPU
Serial GPU

# Quantum Matcha Tea workflow

Quantum circuit
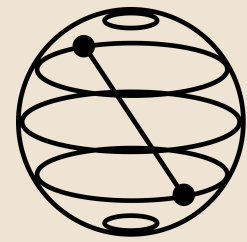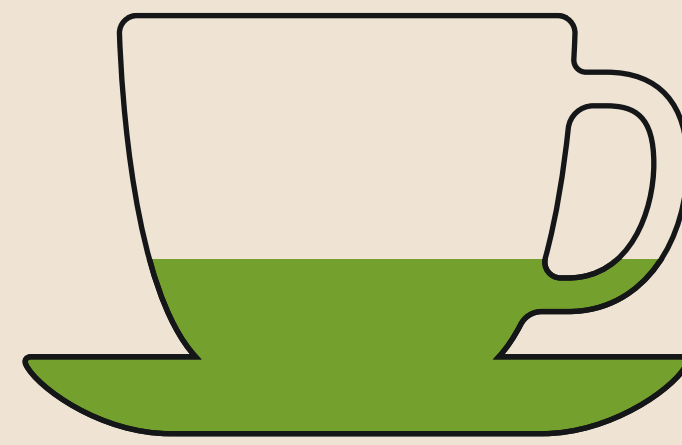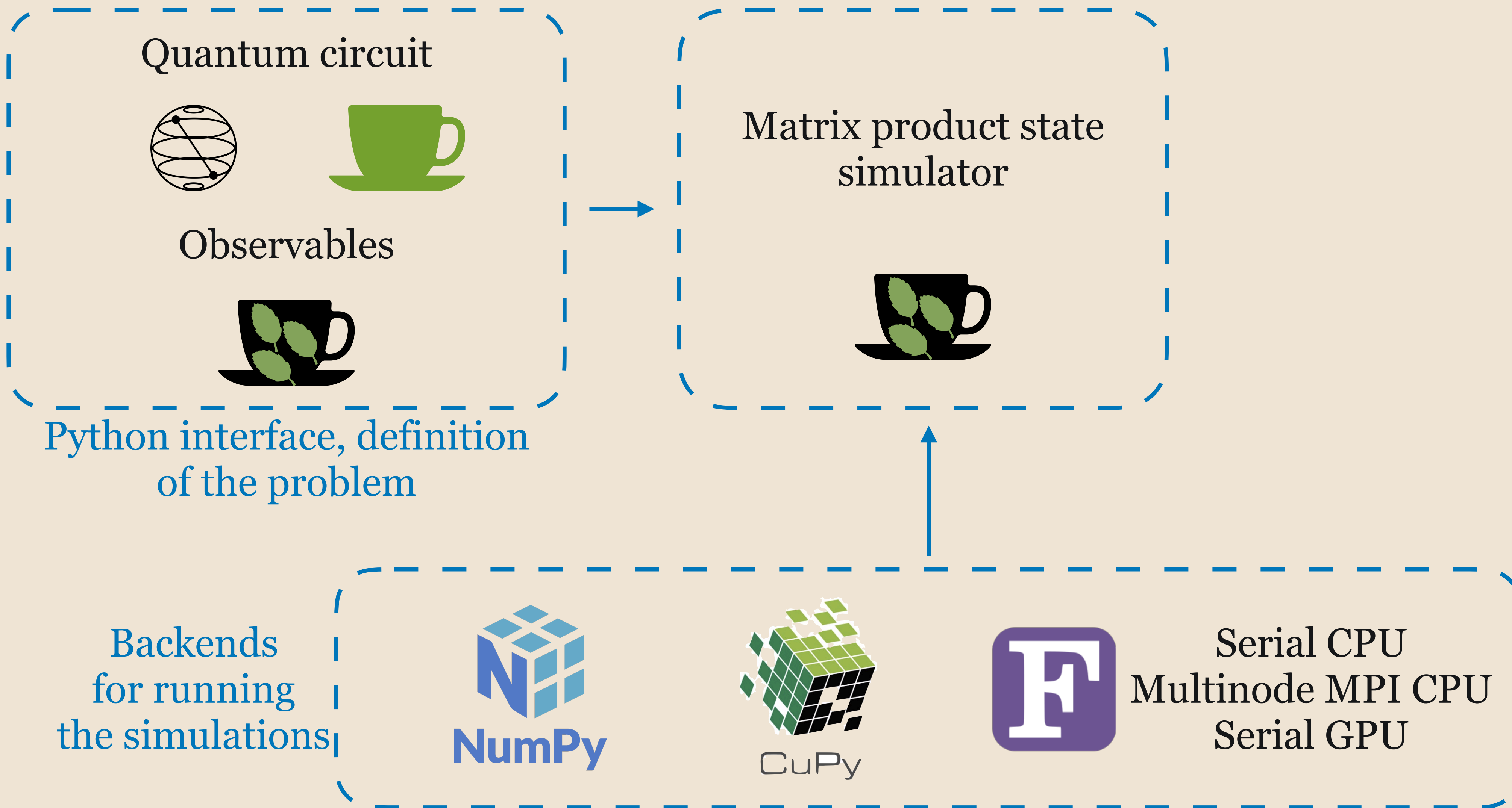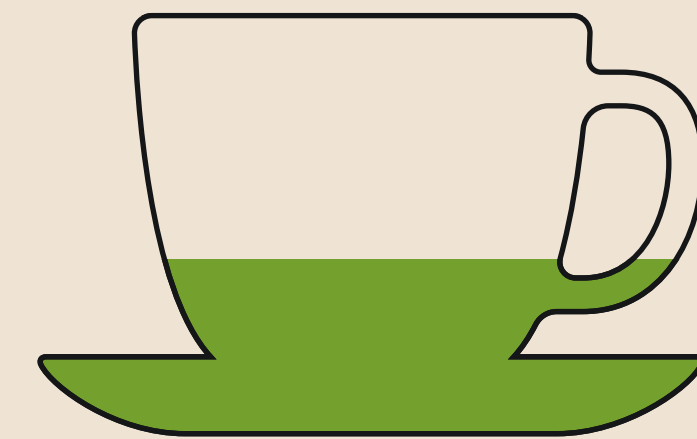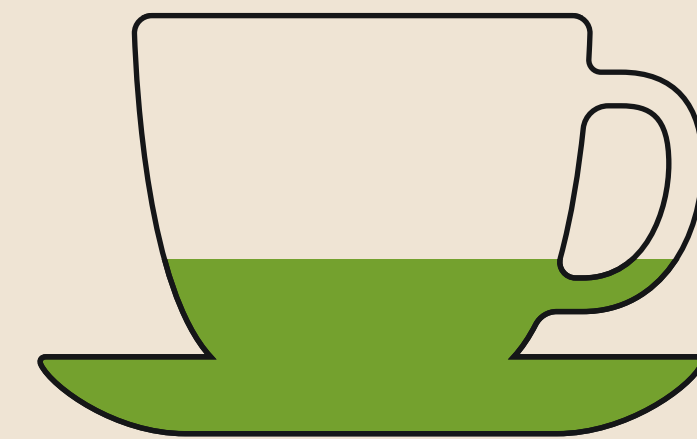
Observables

Matrix product state simulator

Python interface, definition
of the problem

Backends
for running
the simulations

NumPy

CuPy

**F**
Serial CPU
Multinode MPI CPU
Serial GPU

Not public
yet

# Quantum Matcha Tea workflow



Quantum circuit

Observables

Python interface, definition of the problem

Matrix product state simulator

Observables
Runtime statistics
Convergence checks

Python interface output

Backends for running the simulations

NumPy

CuPy

**F** Serial CPU
Multinode MPI CPU
Serial GPU

Not public yet

6

# Convergence checks & error bound

$$|\psi\rangle = \sum_{\alpha=1}^{\chi_T^{i-1}} \boxed{|A_\alpha\rangle} \overset{\lambda_\alpha}{-} \boxed{|B_\alpha\rangle}$$

# Convergence checks & error bound

$$|\psi\rangle = \sum_{\alpha=1}^{\chi_T^{i-1}}$$
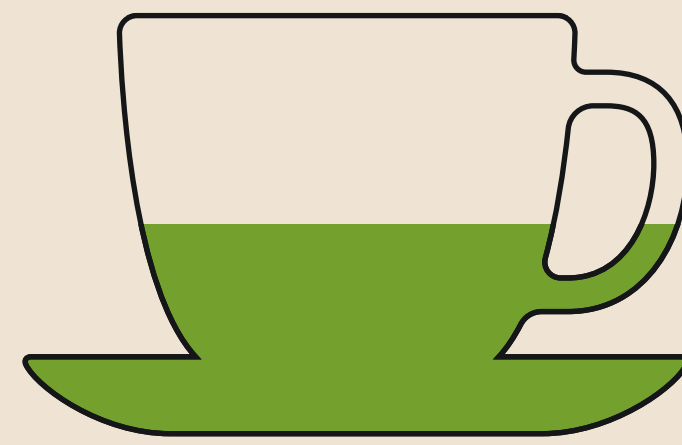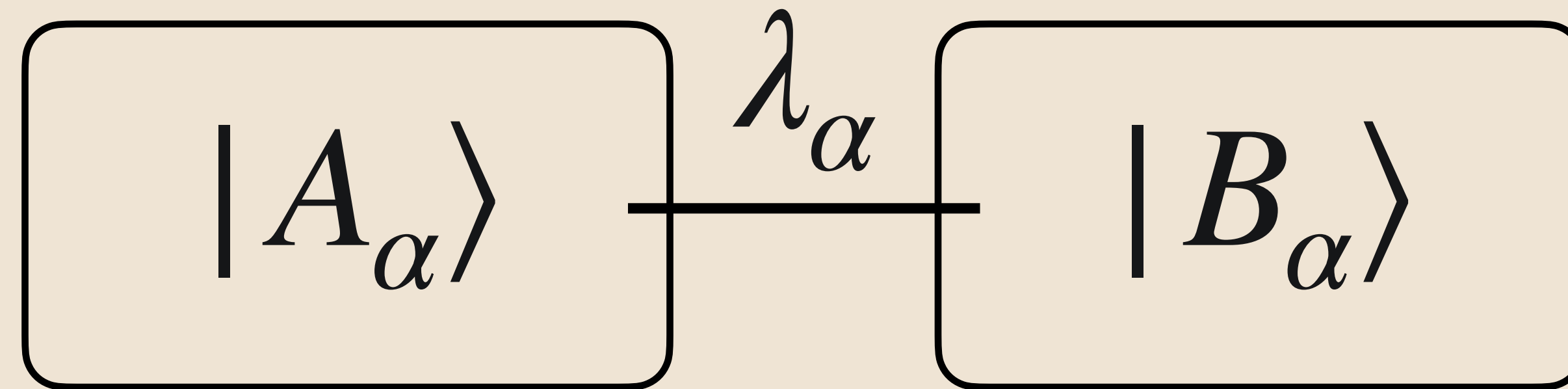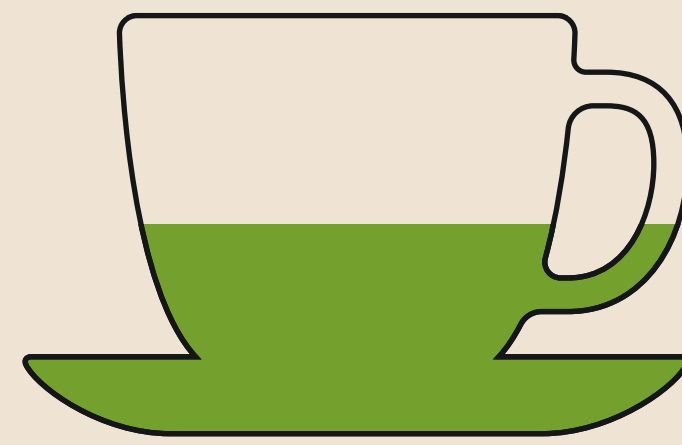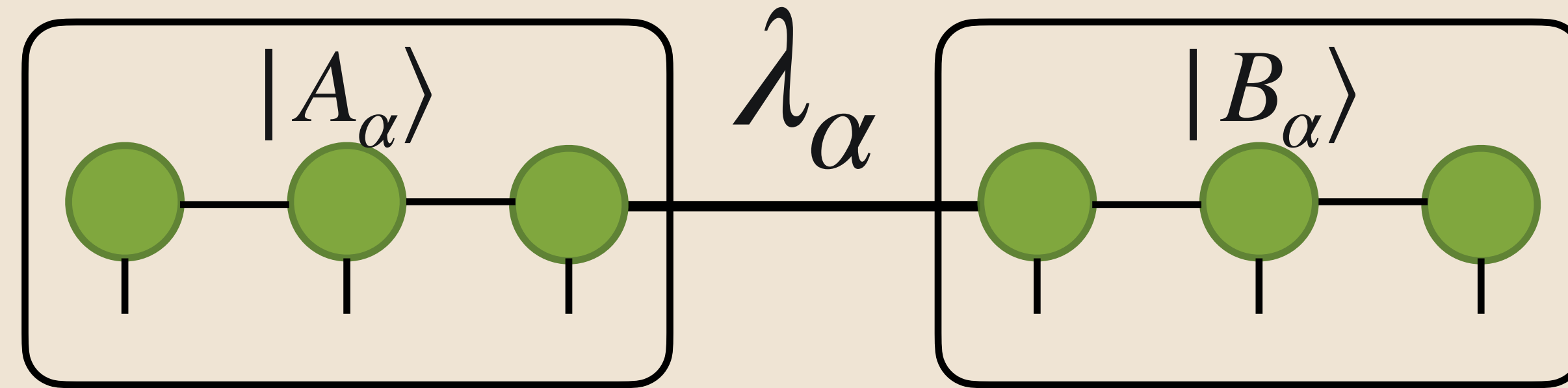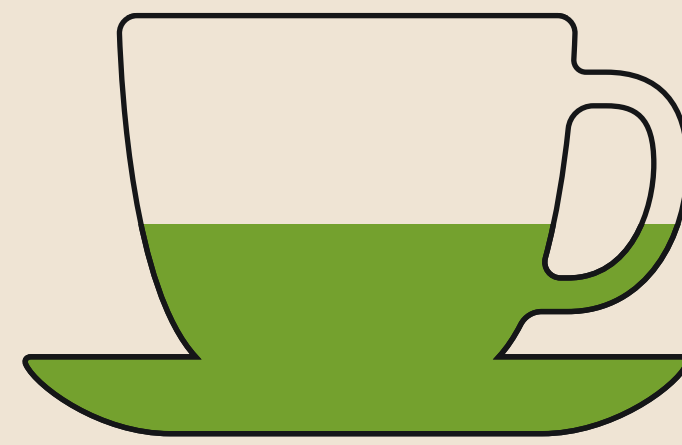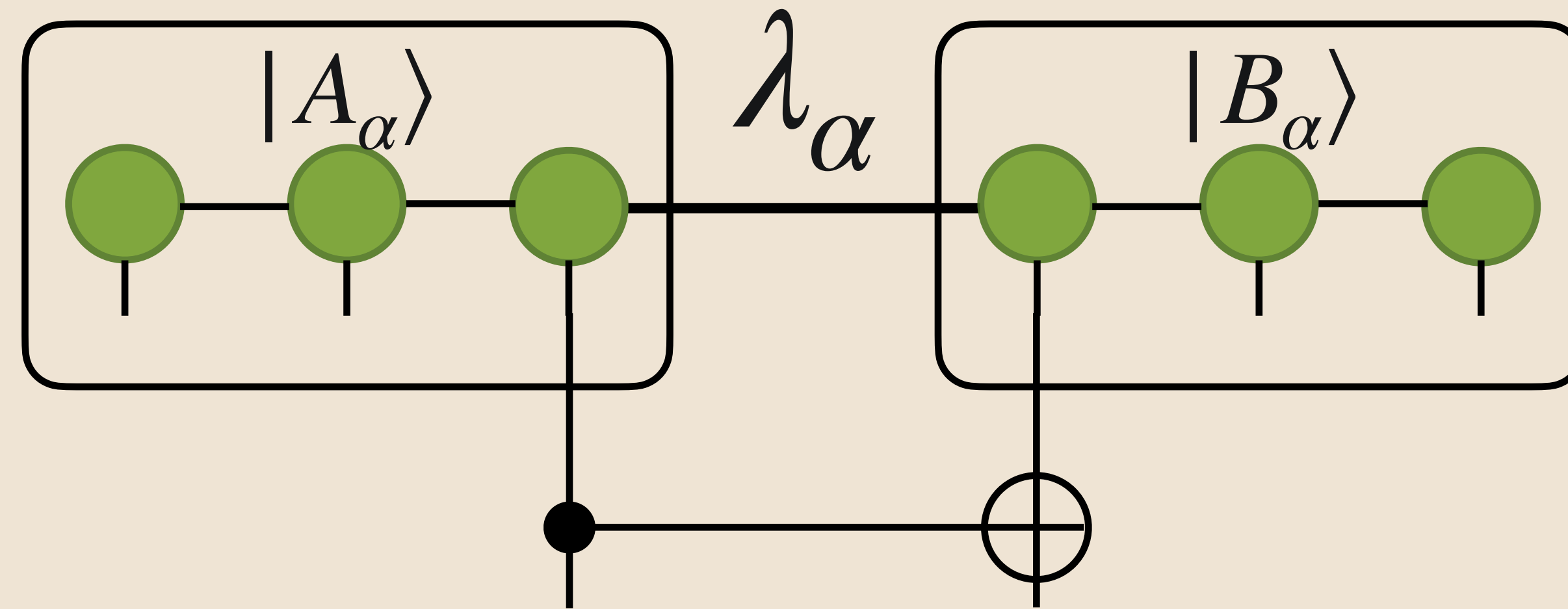
# Convergence checks & error bound

$$|\psi\rangle = \sum_{\alpha=1}^{\chi_T^{i-1}}$$

$$|\psi\rangle = \sum_{\alpha=1}^{\chi_T^{i-1}}$$

Only keep highest $\chi$ singular values, $|\phi\rangle$

# Convergence checks & error bound

$$|\psi\rangle = \sum_{\alpha=1}^{\chi_T^i} \quad \boxed{\substack{|A_\alpha\rangle \\ \bullet - \bullet - \bullet}} \quad \lambda_\alpha \quad \boxed{\substack{|B_\alpha\rangle \\ \bullet - \bullet - \bullet}}$$
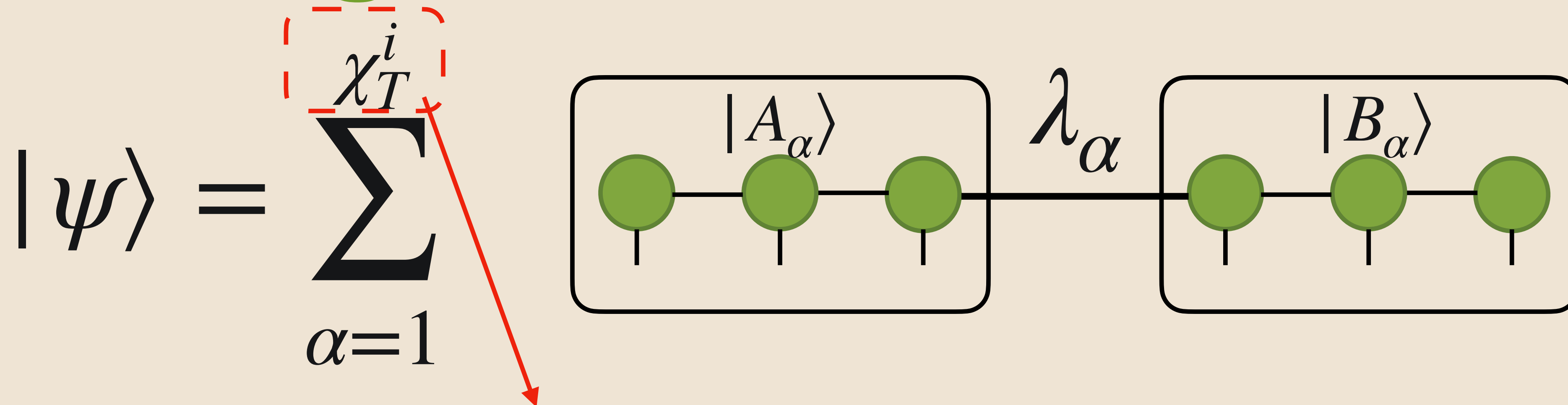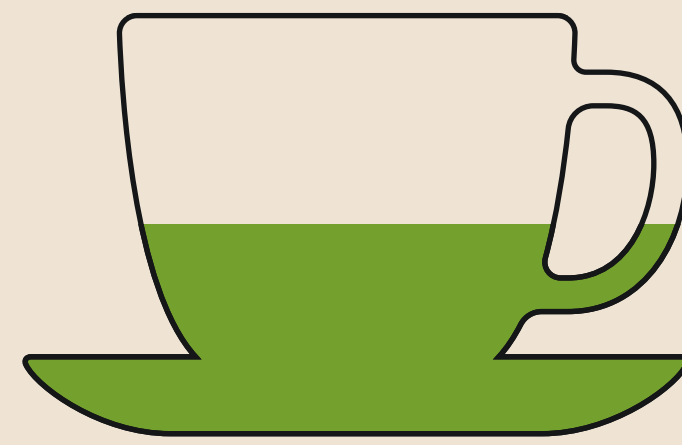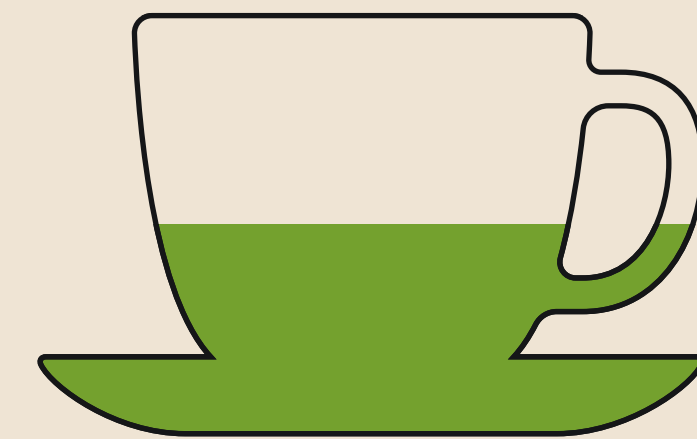
Only keep highest $\chi$ singular values, $|\phi\rangle$

# Convergence checks & error bound

$$|\psi\rangle = \sum_{\alpha=1}^{\chi_T^i} \quad$$



$|A_\alpha\rangle$    $\lambda_\alpha$    $|B_\alpha\rangle$

Only keep highest $\chi$ singular values, $|\phi\rangle$

Fidelity of the state

$$\mathscr{F}_i(\chi) = \left| \langle \psi | \phi \rangle \right|^2 = \left| 1 - \sum_{\alpha=\chi+1}^{\chi_T^i} \lambda_\alpha^2 \right|^2$$

# Convergence checks & error bound

$$|\psi\rangle = \sum_{\alpha=1}^{\chi_T^i}$$



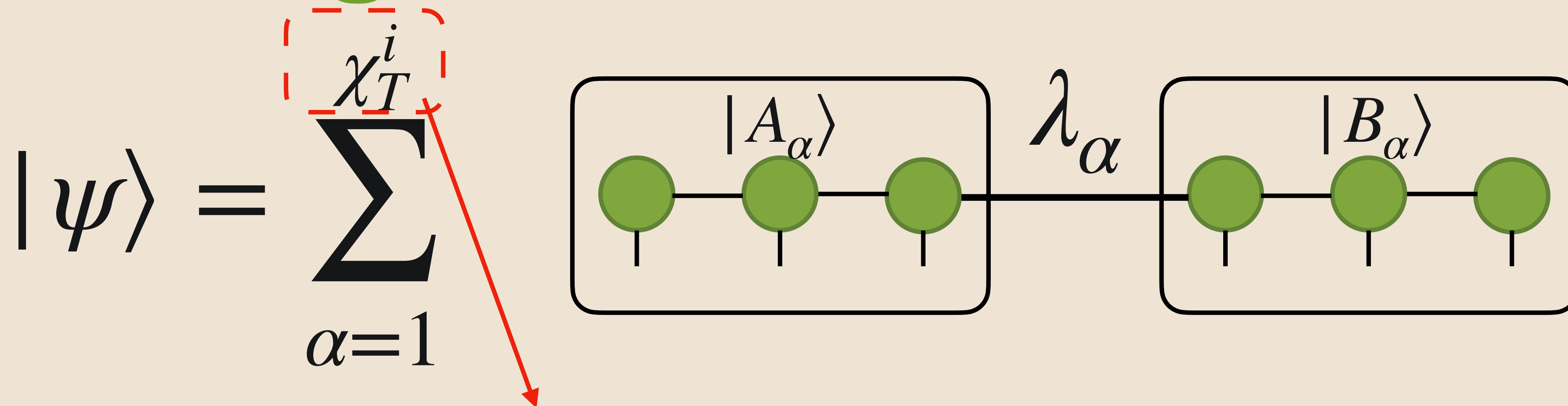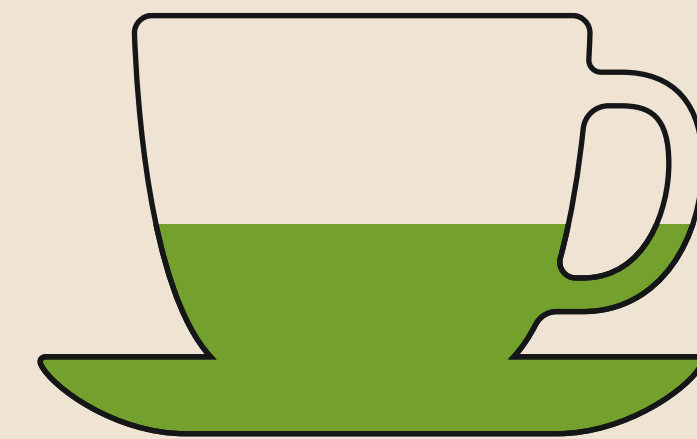$|A_\alpha\rangle$   $\lambda_\alpha$   $|B_\alpha\rangle$

Only keep highest $\chi$ singular values, $|\phi\rangle$

Fidelity of the state
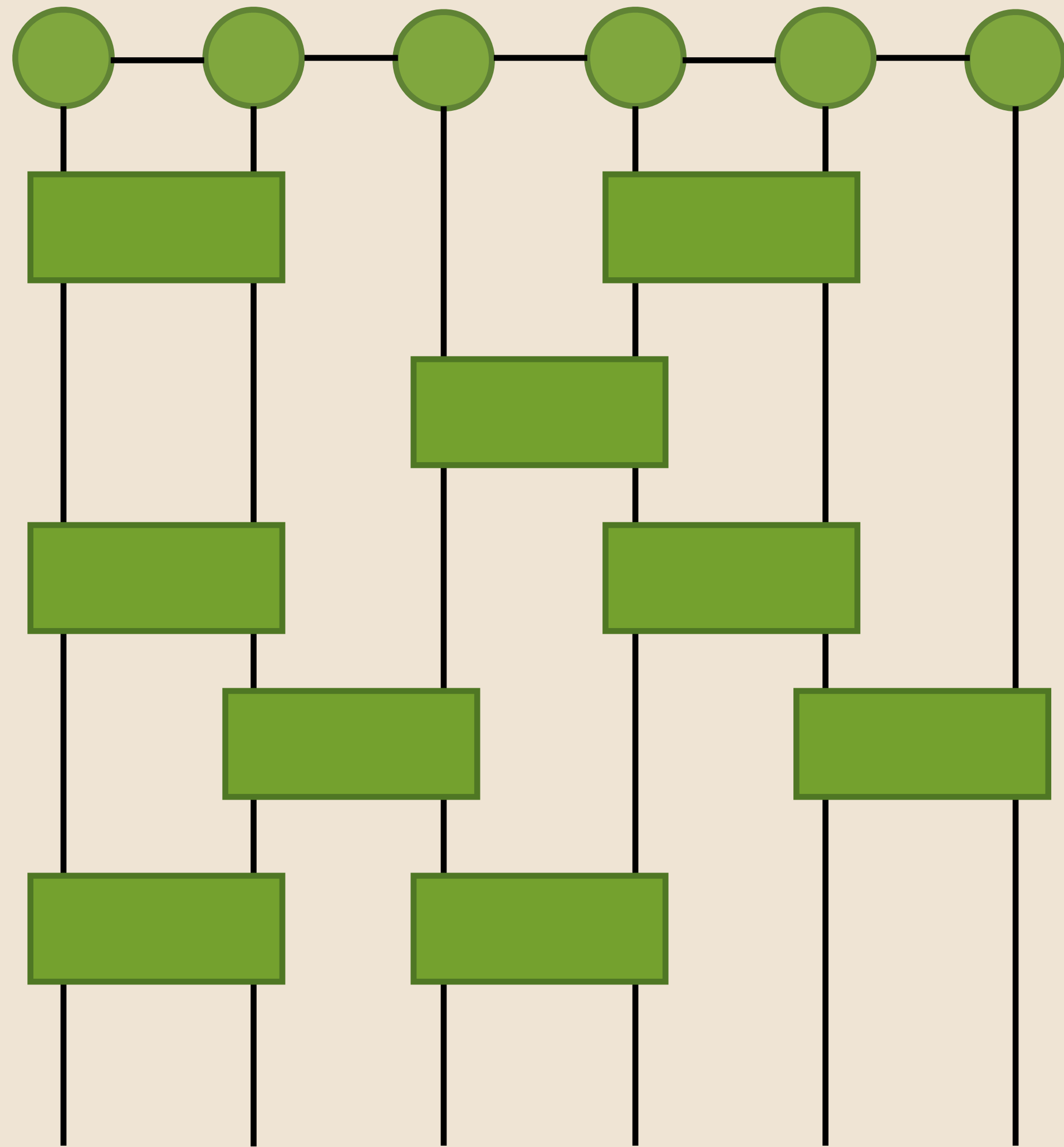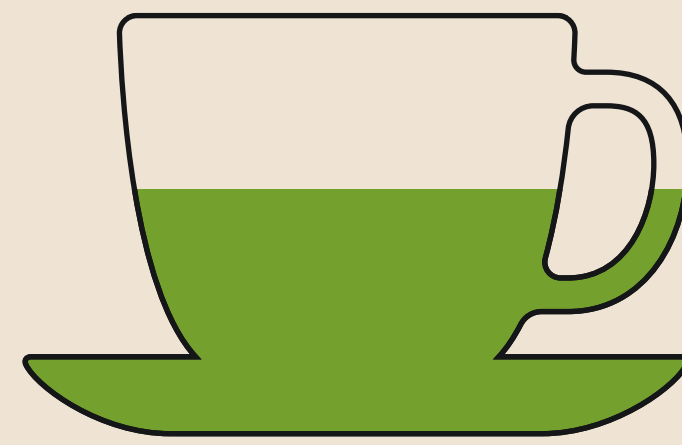
$$\mathscr{F}_i(\chi) = \left| \langle\psi|\phi\rangle \right|^2 = \left| 1 - \sum_{\alpha=\chi+1}^{\chi_T^i} \lambda_\alpha^2 \right|^2$$
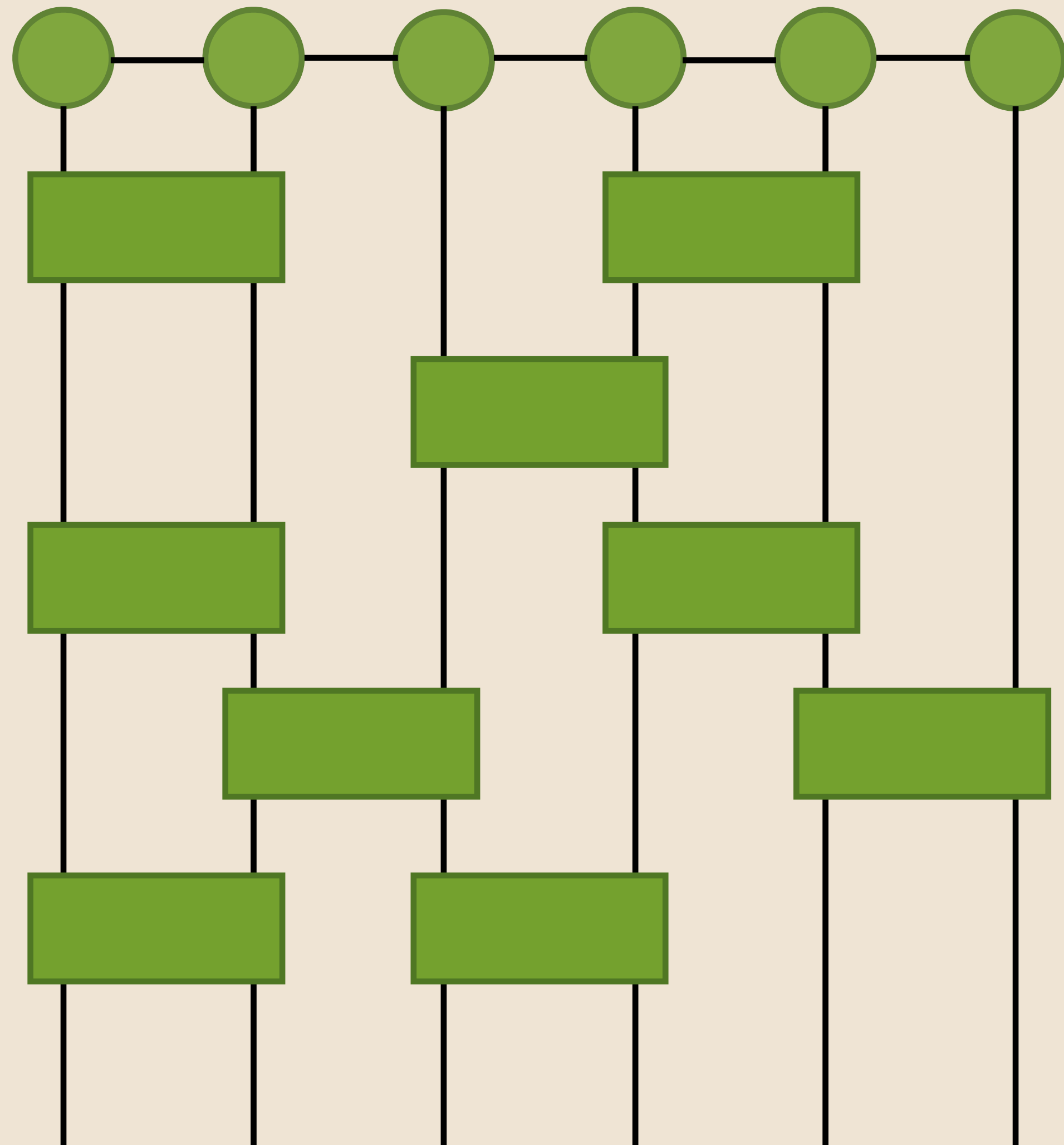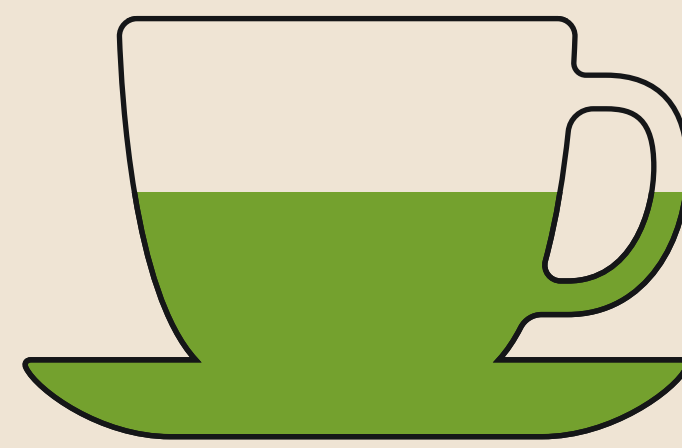
Computed during the simulation

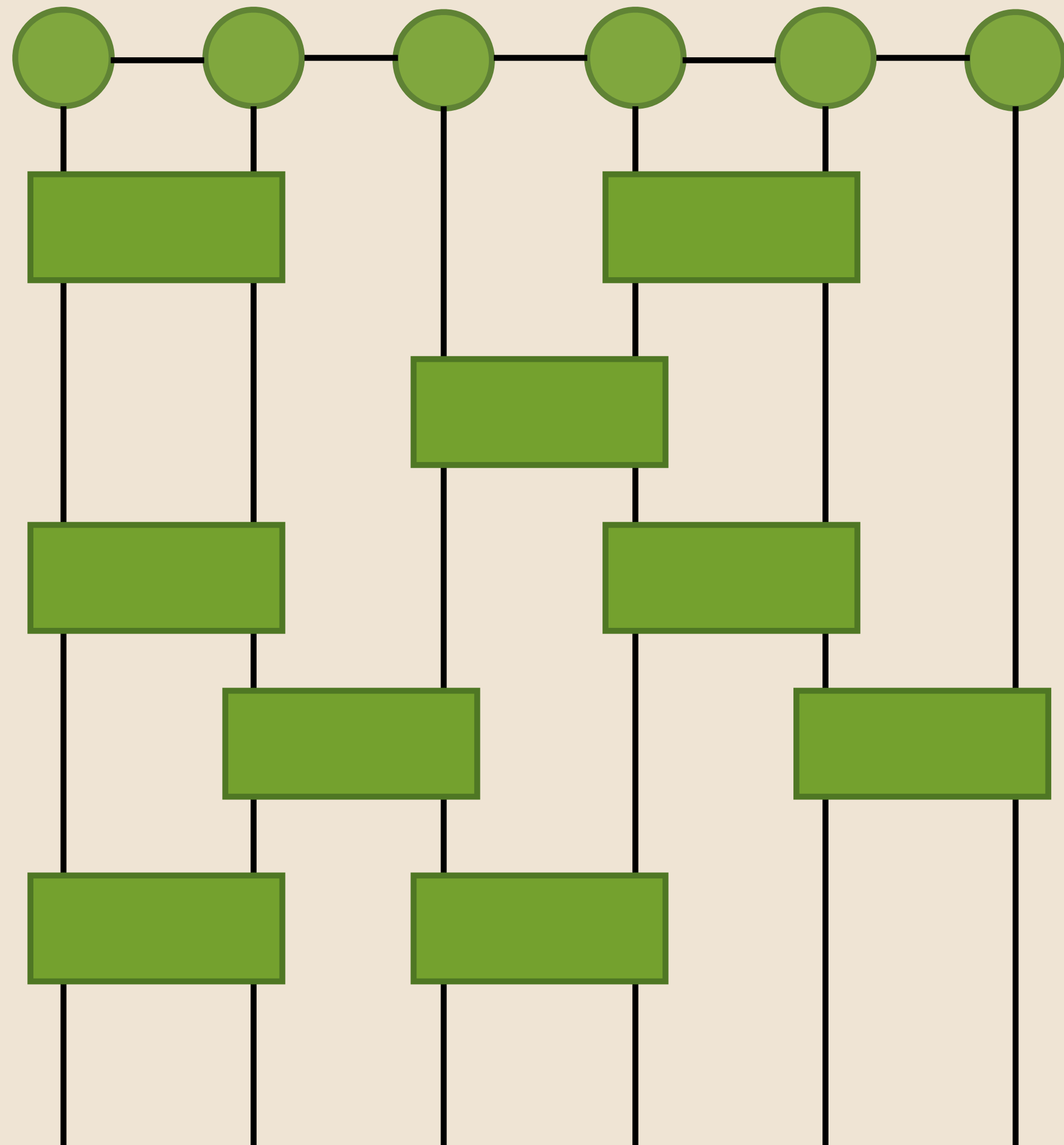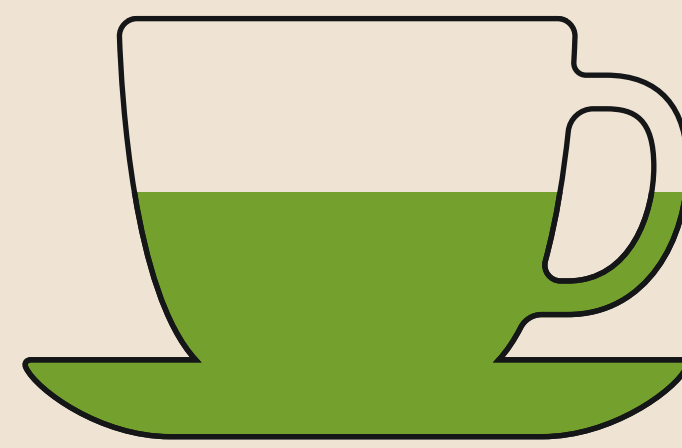# Convergence and error checks

# Convergence and error checks



Fidelity of the state after a **single** gate

$$\mathscr{F}_i(\chi) = \left| 1 - \sum_{\alpha = {\color{red}\chi + 1}}^{\chi_T^i} \lambda_\alpha^2 \right|^2$$

# Convergence and error checks
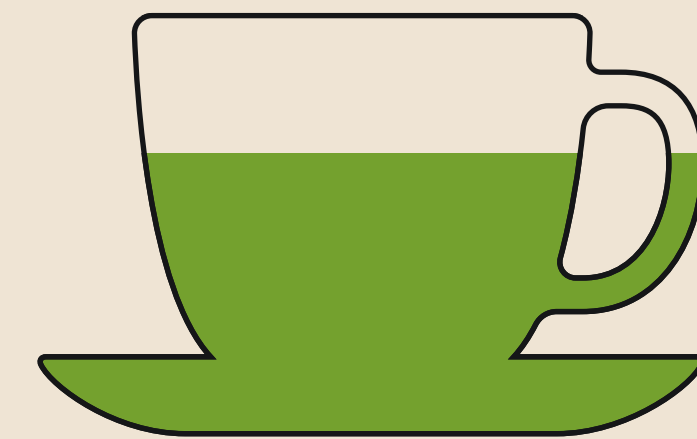
Fidelity of the state after a **single** gate

$$\mathcal{F}_i(\chi) = \left| 1 - \sum_{\alpha=\chi+1}^{\chi_T^i} \lambda_\alpha^2 \right|^2$$
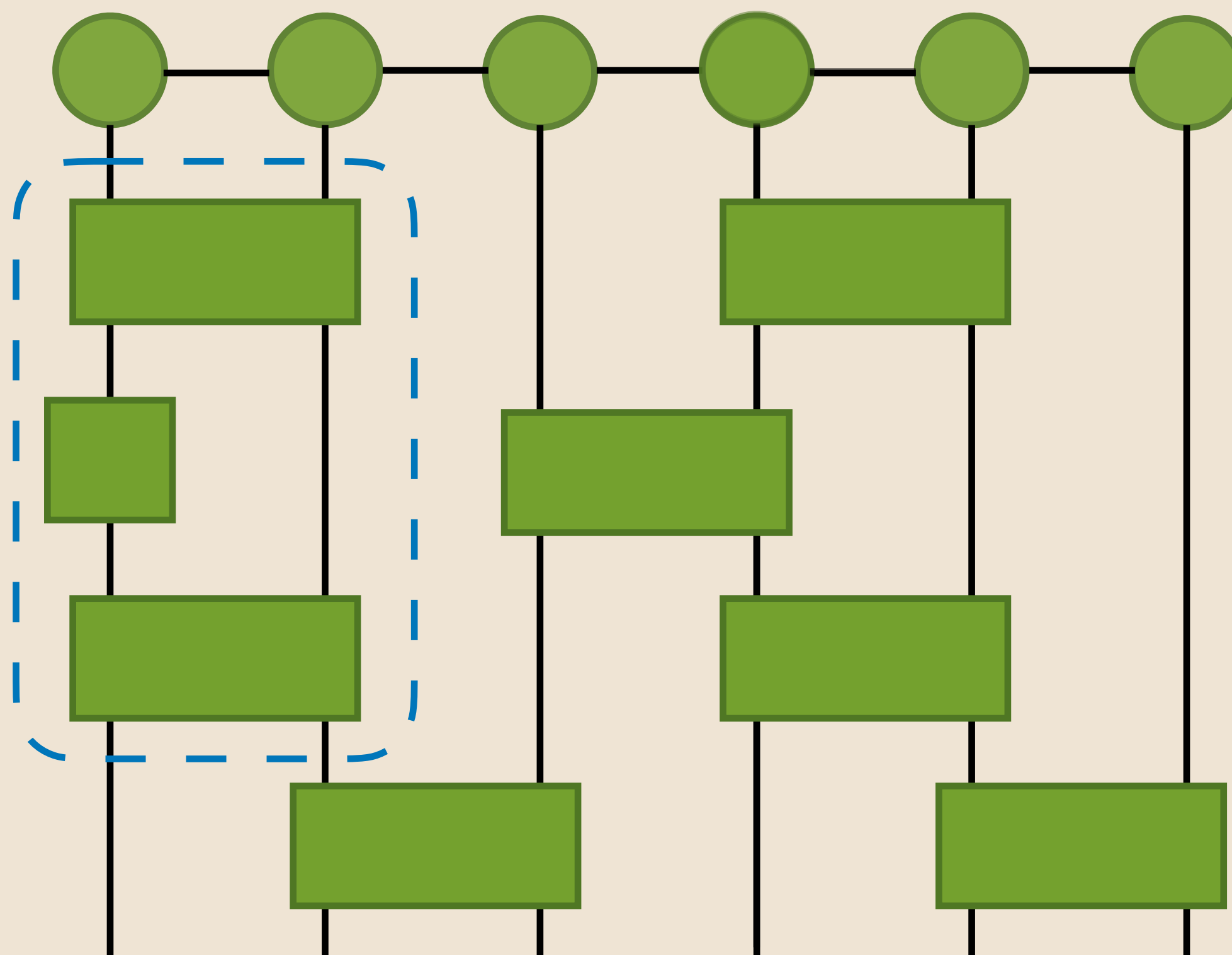
Fidelity at the end of the simulation

$$\mathcal{F}^{tot}(\chi) \geq \prod_i \mathcal{F}_i(\chi)$$

8

# Optimisation & parallelism



Gates acts on the same qubits:
we contract gates together and only after with state

# Optimisation & parallelism



Gates acts on the same qubits:
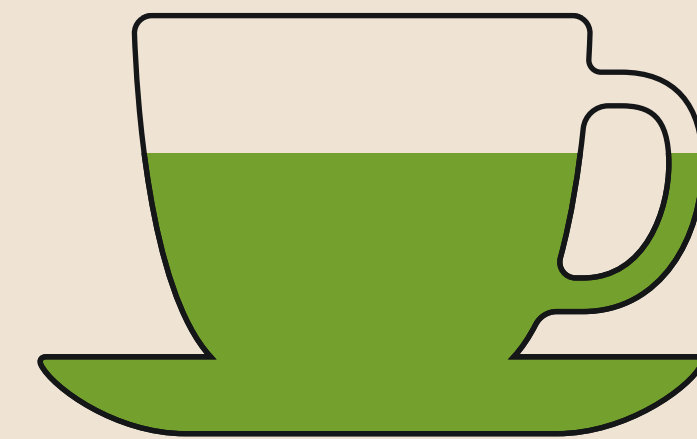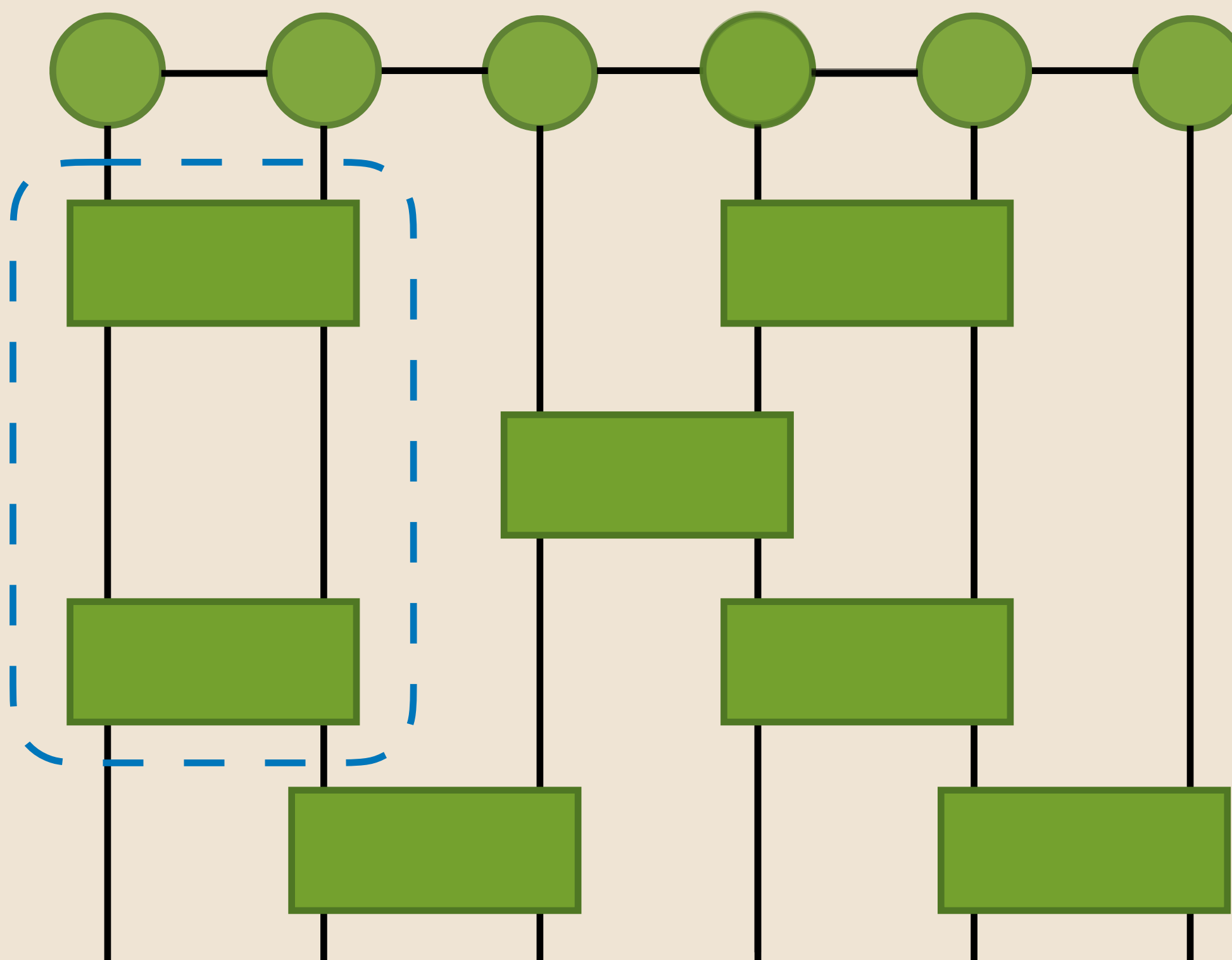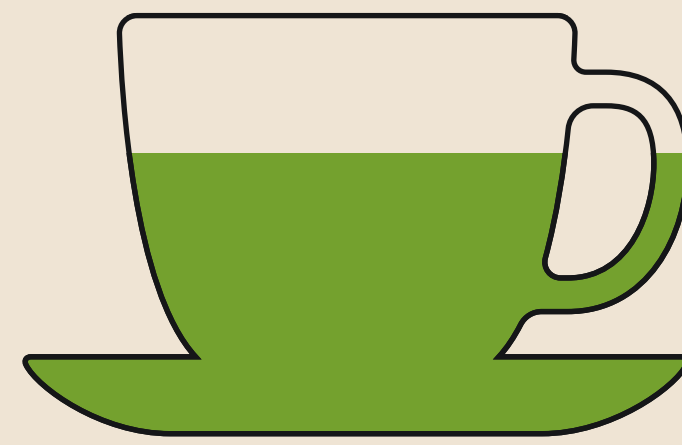we contract gates together and only after with state

# Optimisation & parallelism



Gates acts on the
same qubits:
we contract gates
together and only
after with state

# Optimisation & parallelism



Node 0          Node 1

Gates acts on the
same qubits:
we contract gates
together and only
after with state

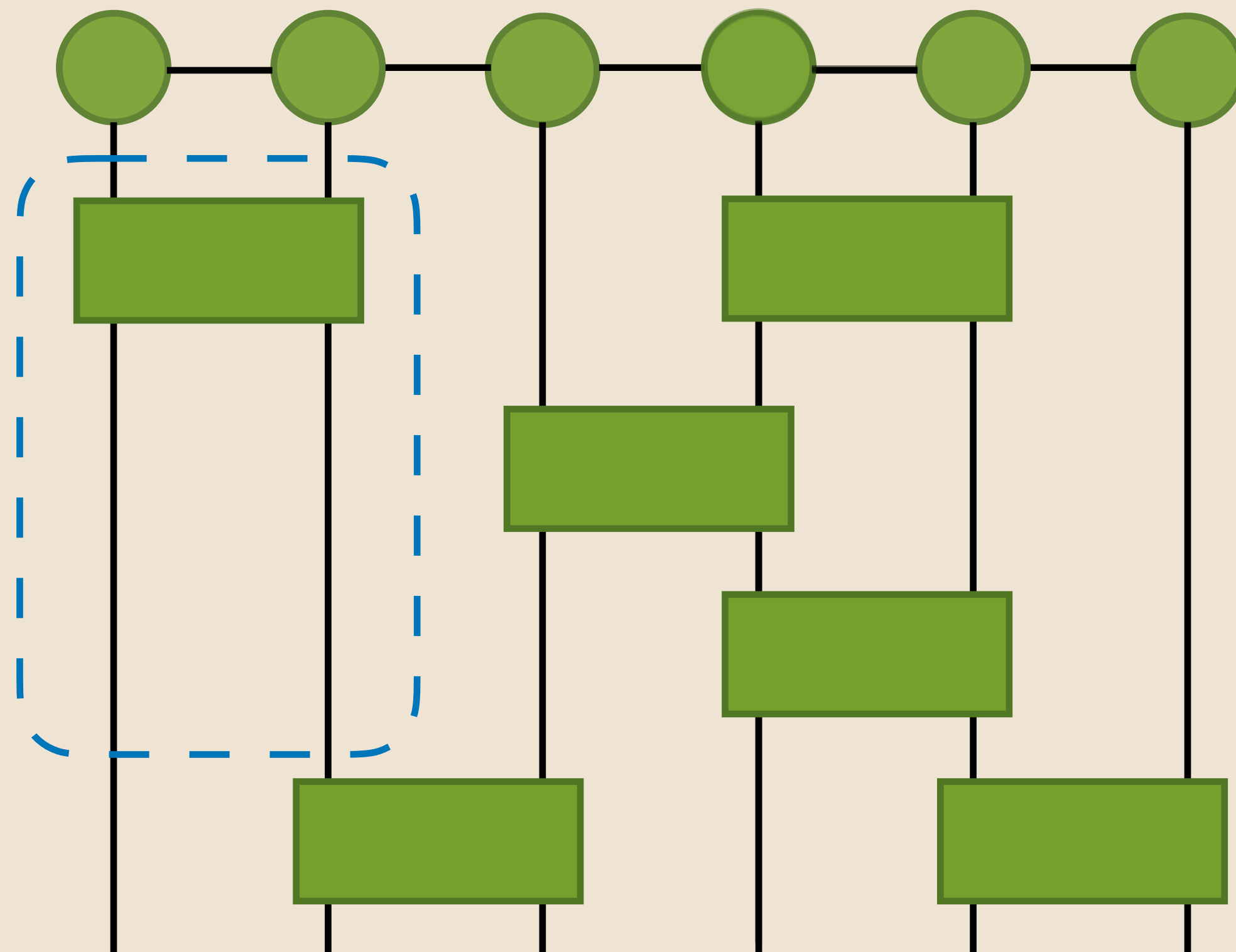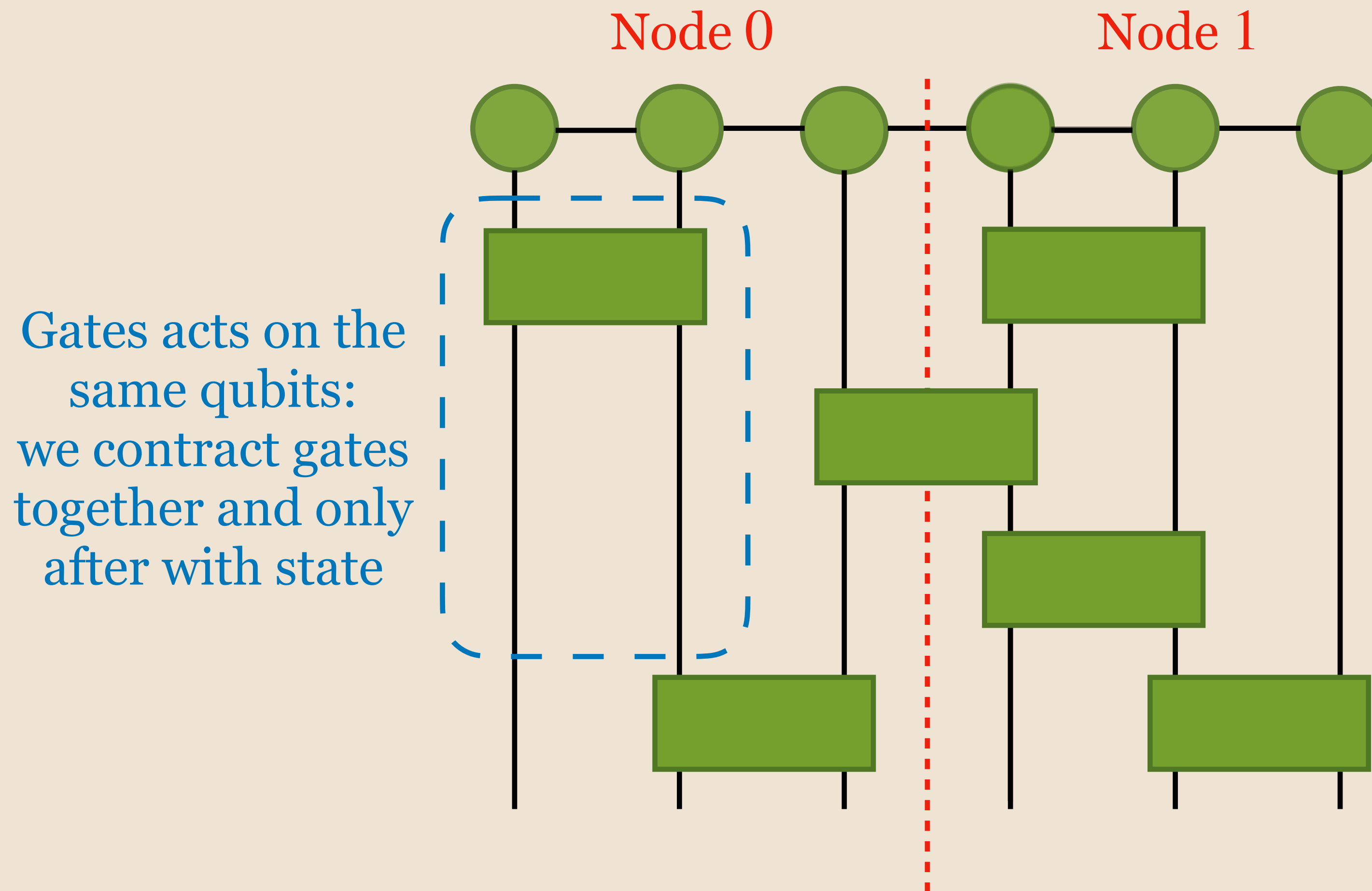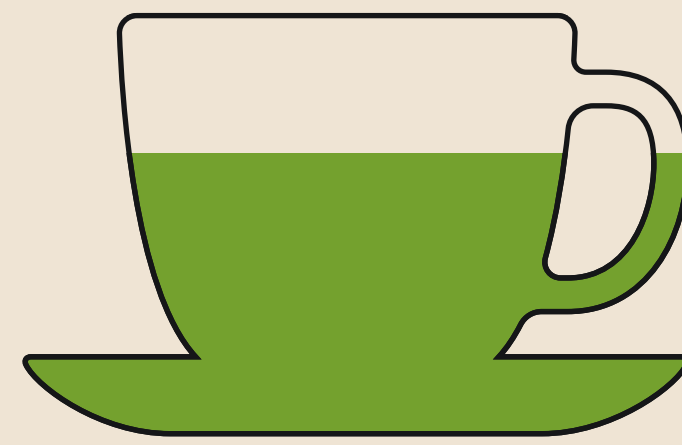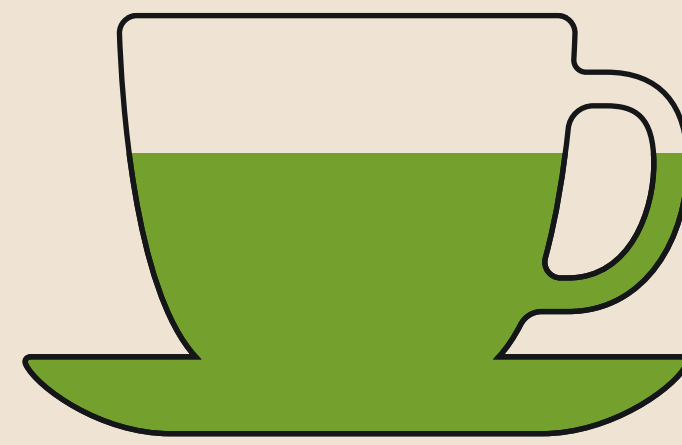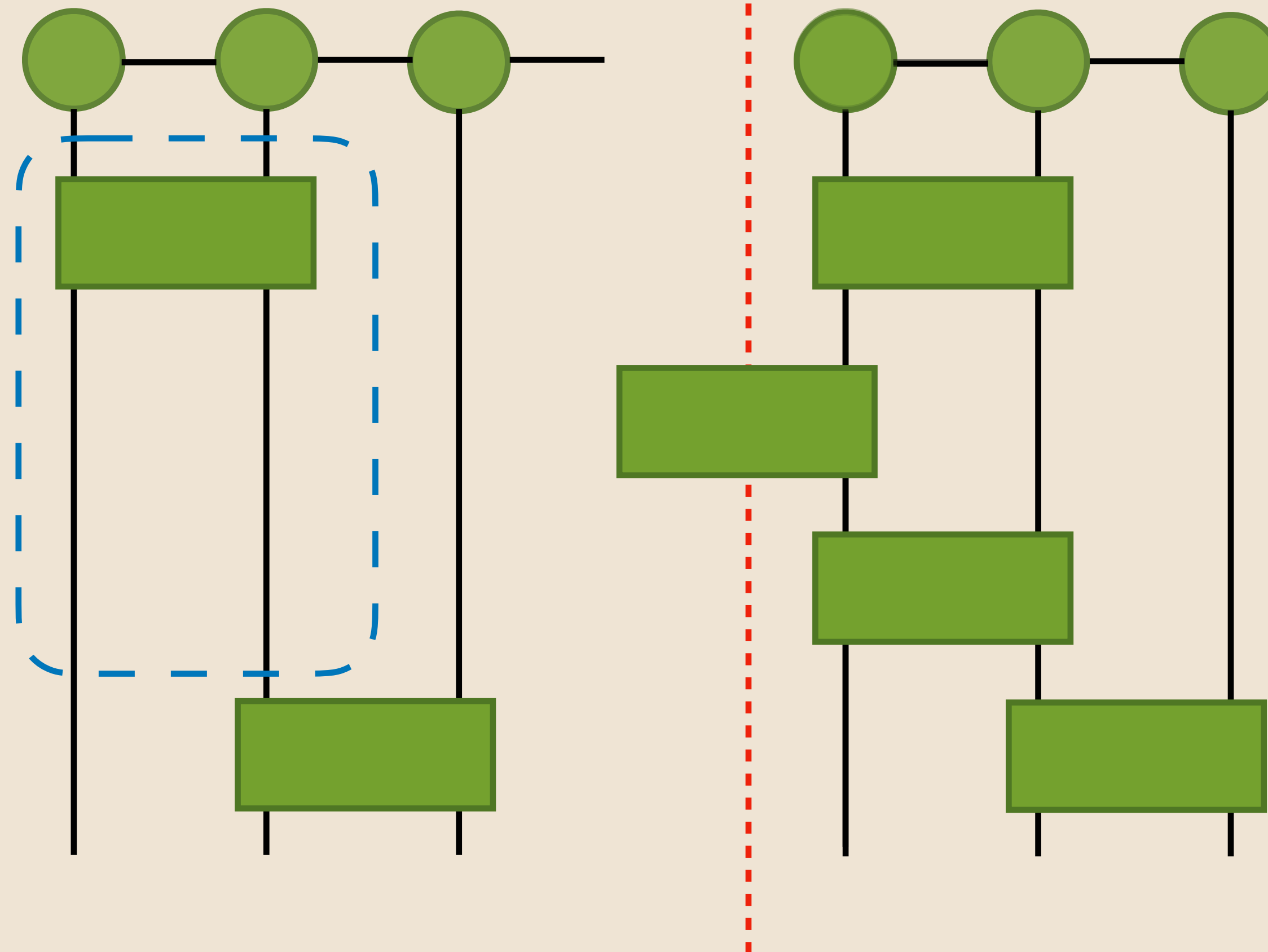# Optimisation & parallelism

Node 0

Node 1

Gates acts on the same qubits:
we contract gates together and only after with state

# Optimisation & parallelism



Copy of the qubit state
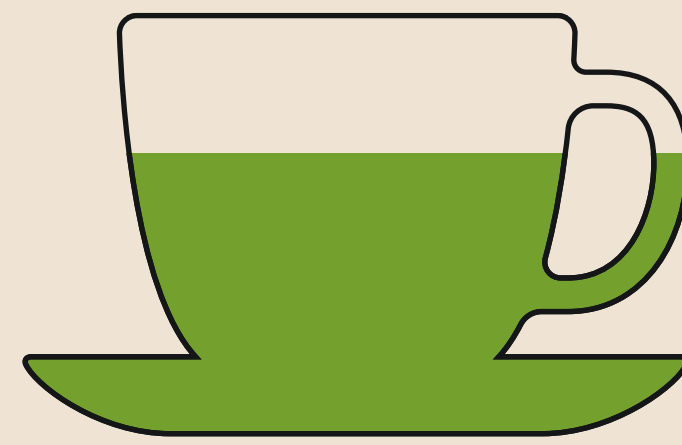
Node 0     Node 1

Gates acts on the
same qubits:
we contract gates
together and only
after with state

# Optimisation & parallelism

Copy of the qubit state

Node 0          Node 1

Gates acts on the
same qubits:
we contract gates
together and only
after with state

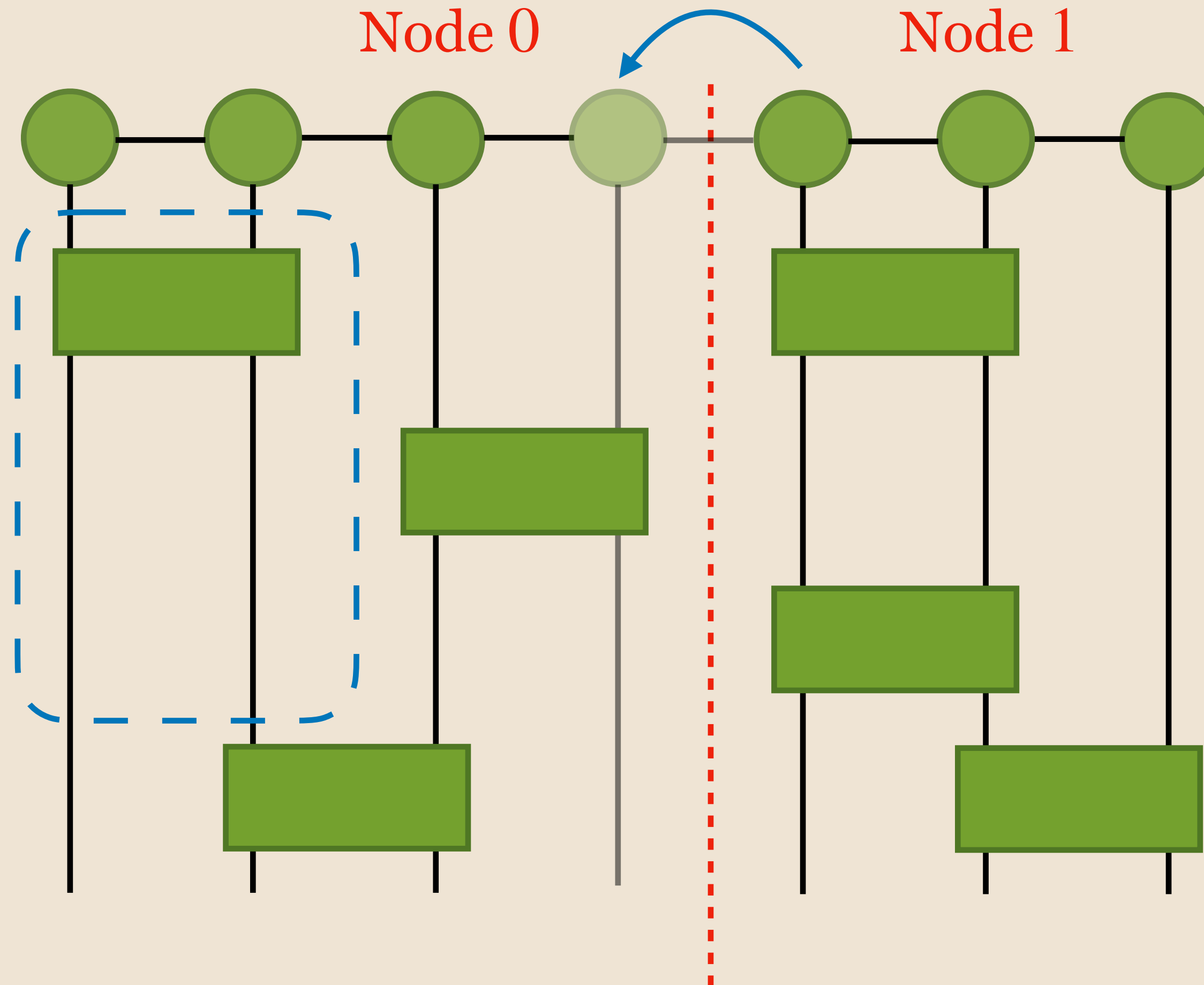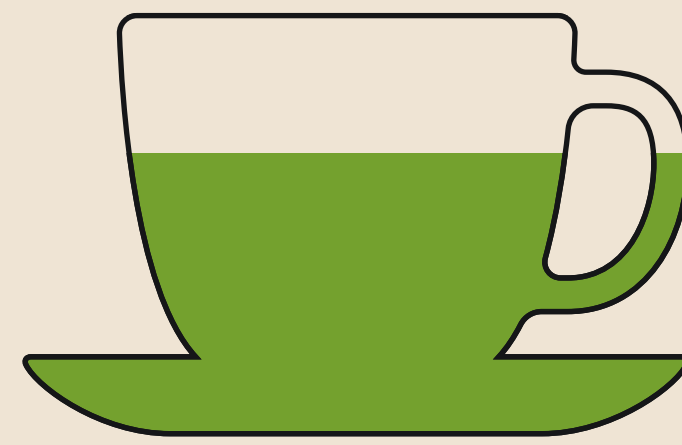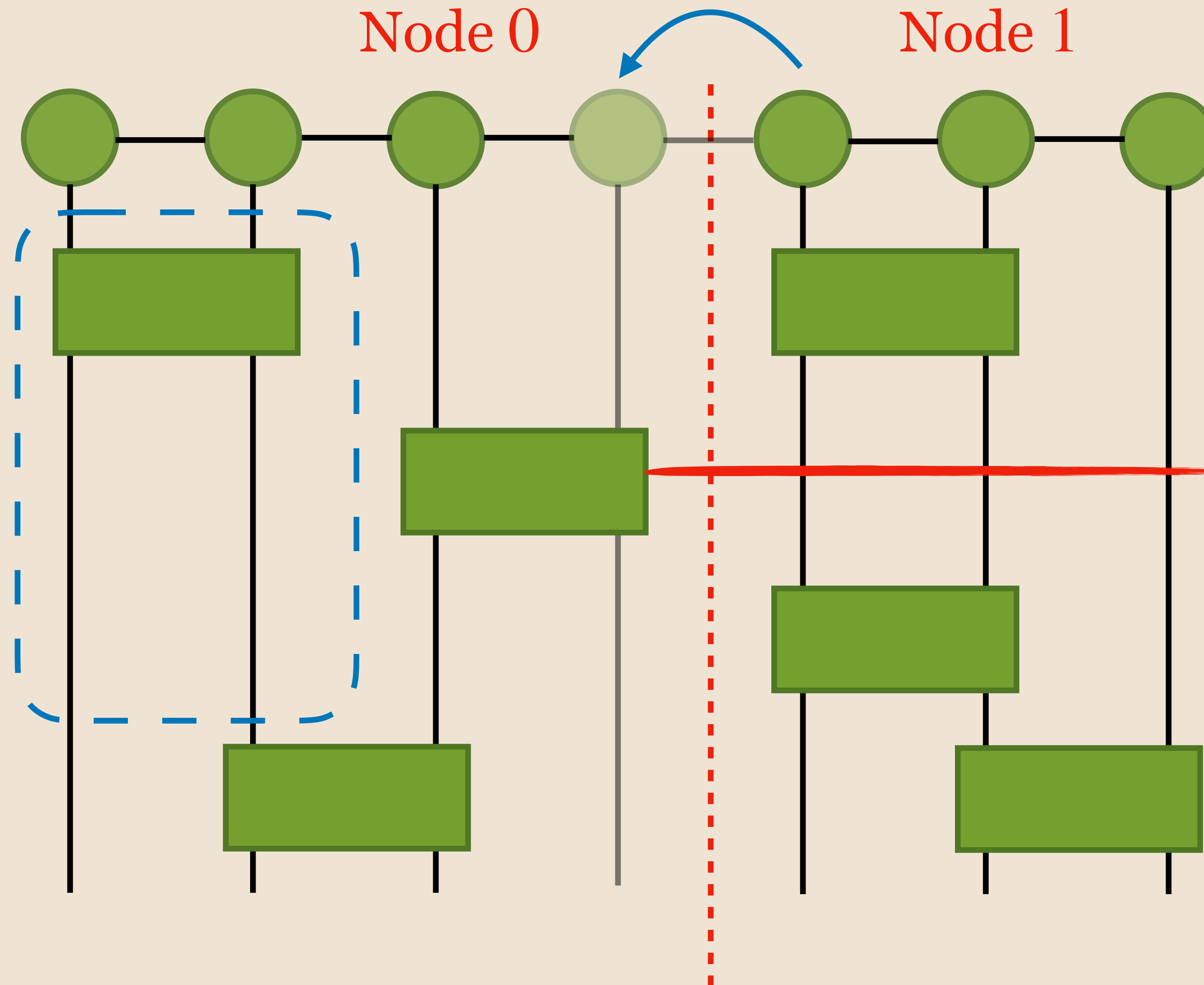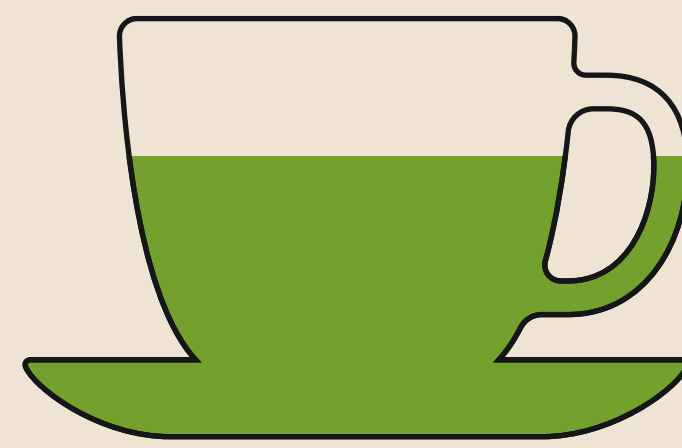Barrier to wait for the
data from node 0

# Optimisation & parallelism



Copy of the qubit state

Node 0

Node 1

Gates acts on the same qubits: we contract gates together and only after with state

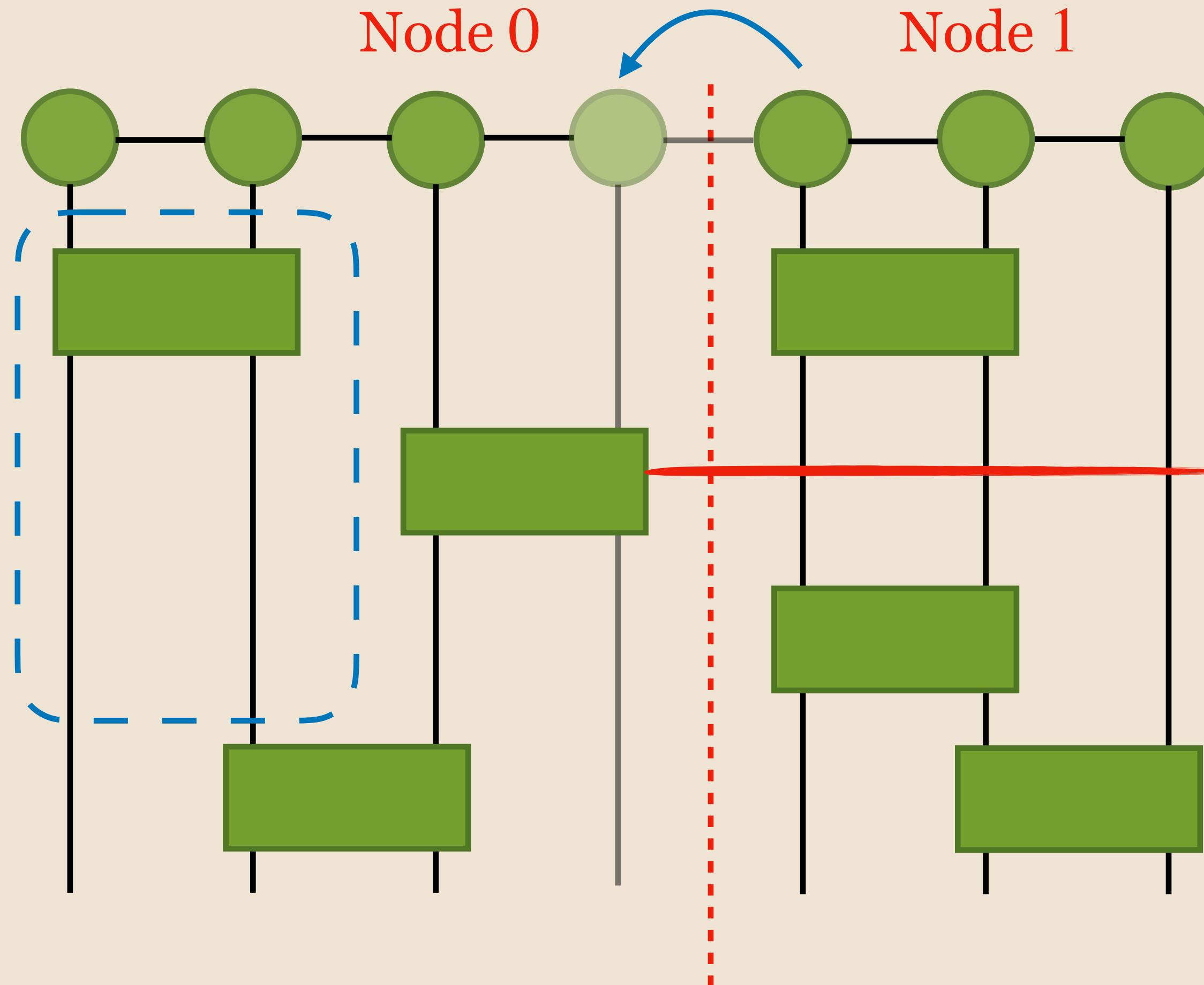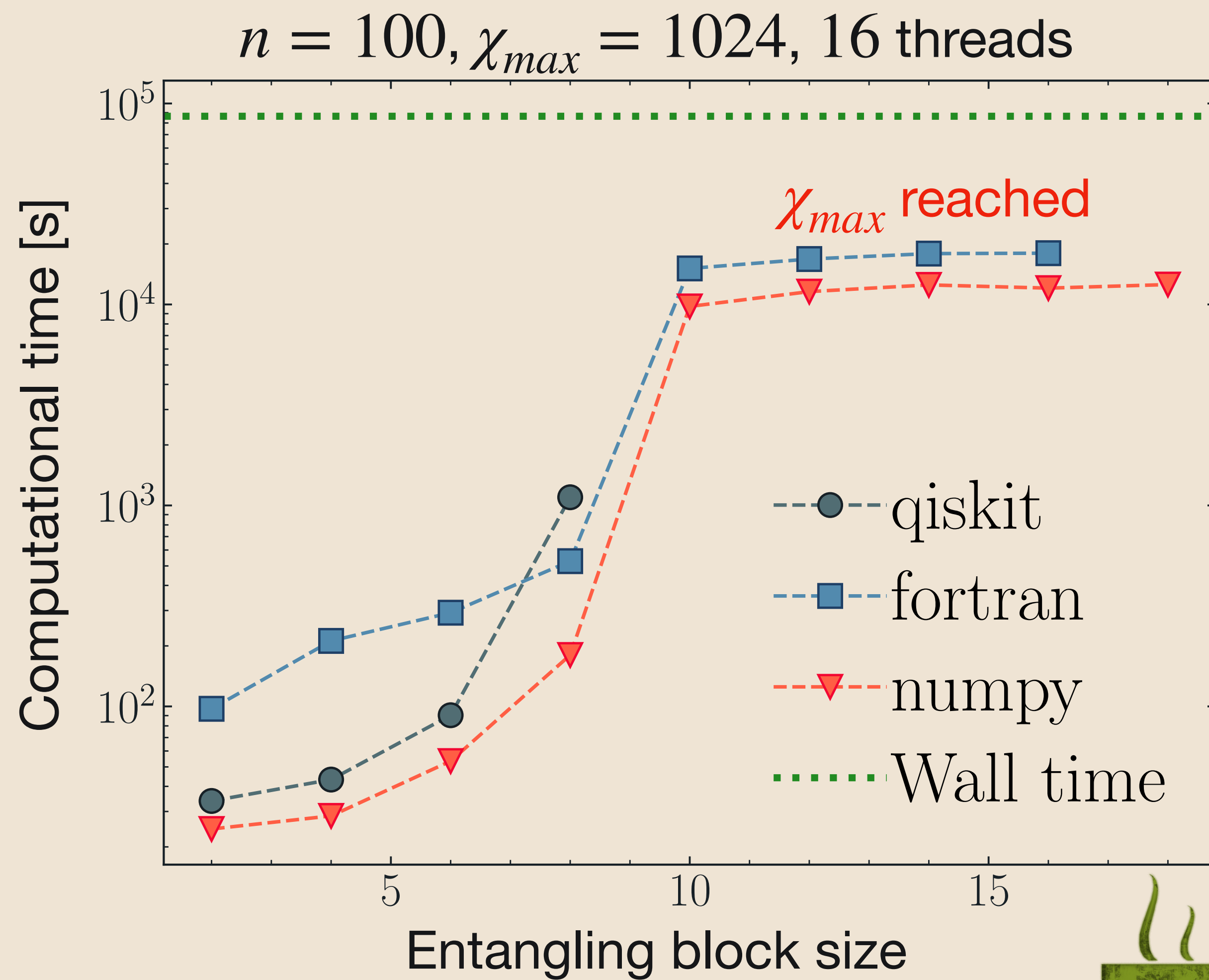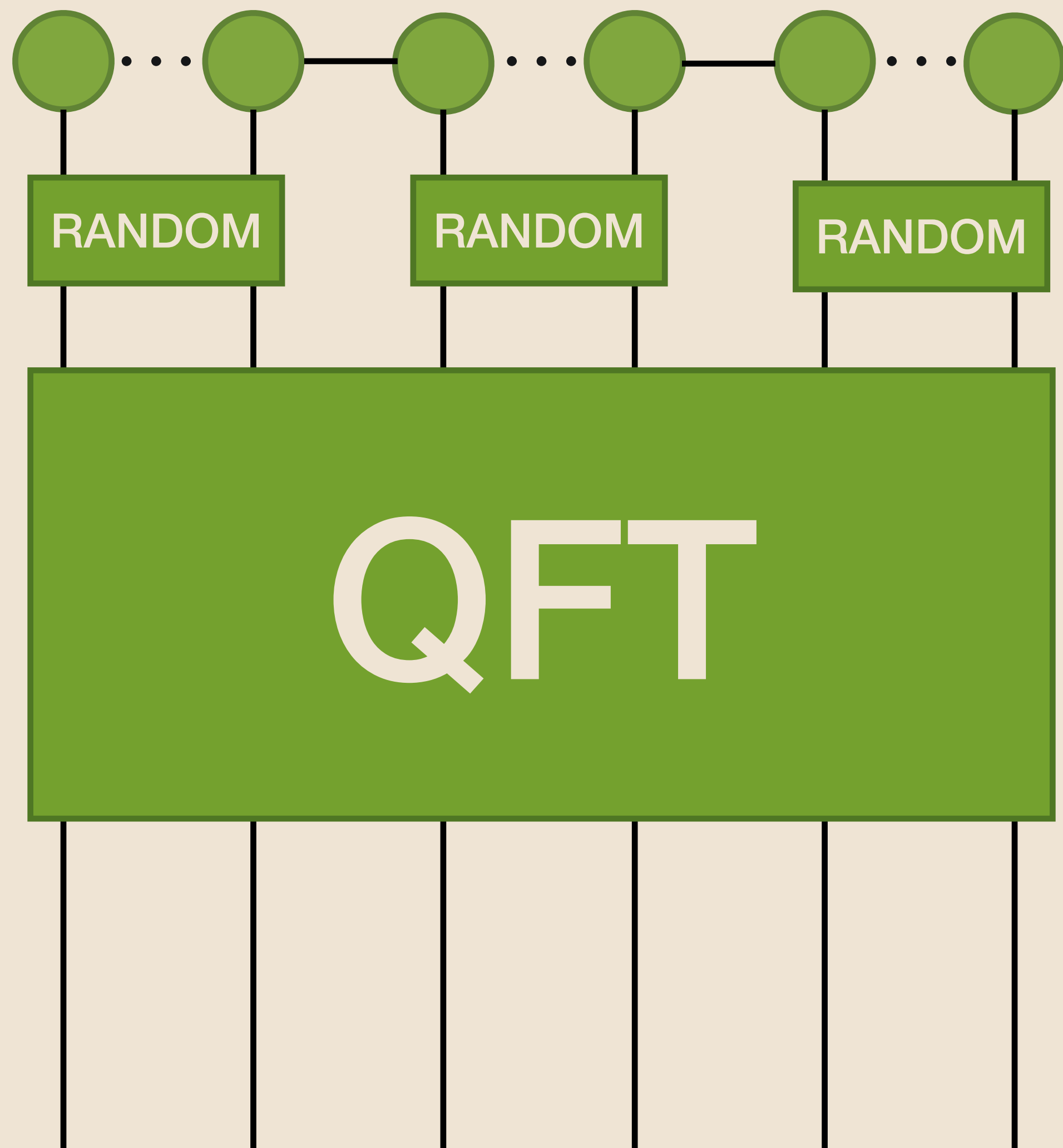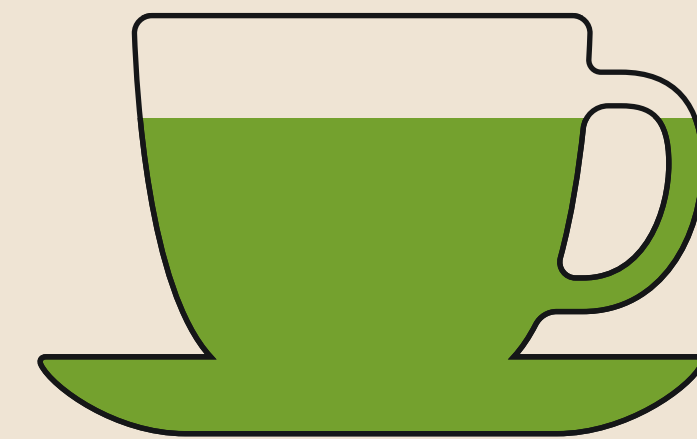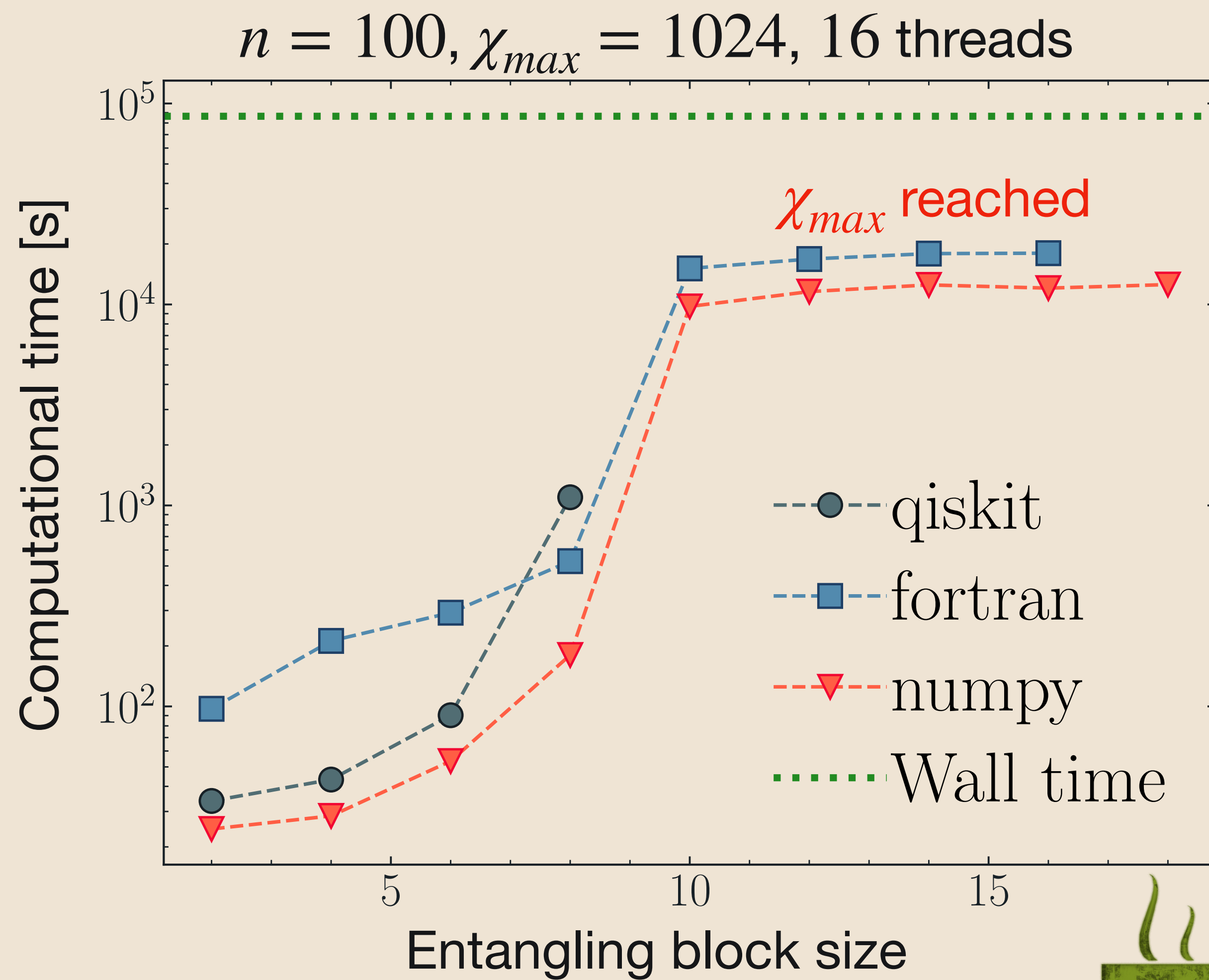Barrier to wait for the data from node 0

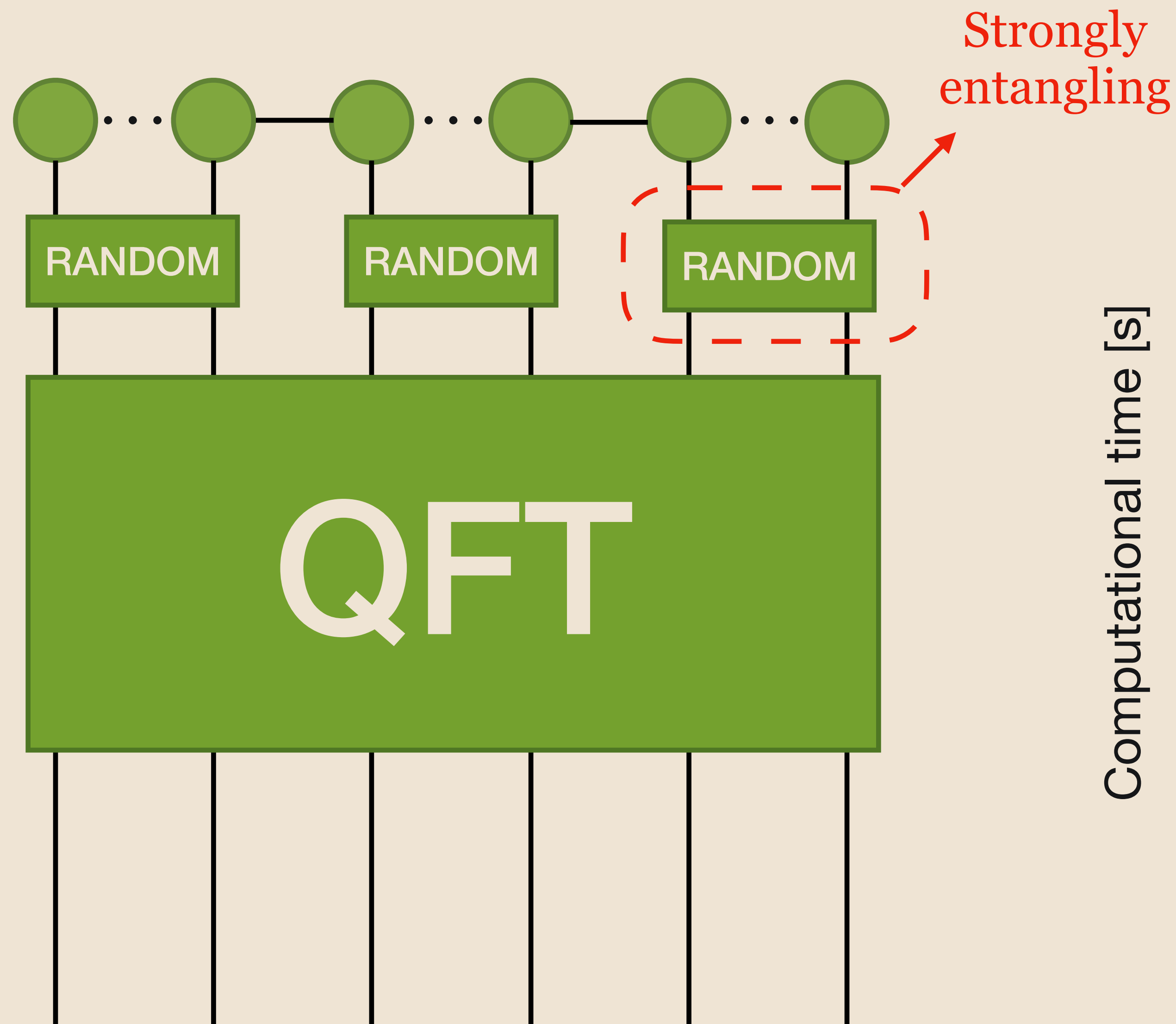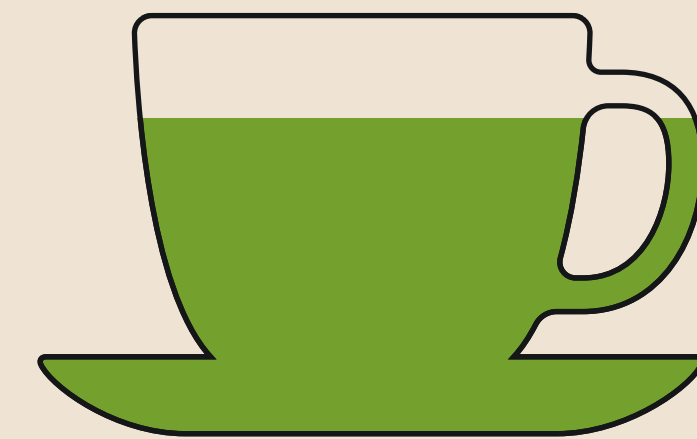A GOOD PARALLEL SCALING INCREASES ERRORS DUE TO AN ALGORITHMIC SUBTLETY

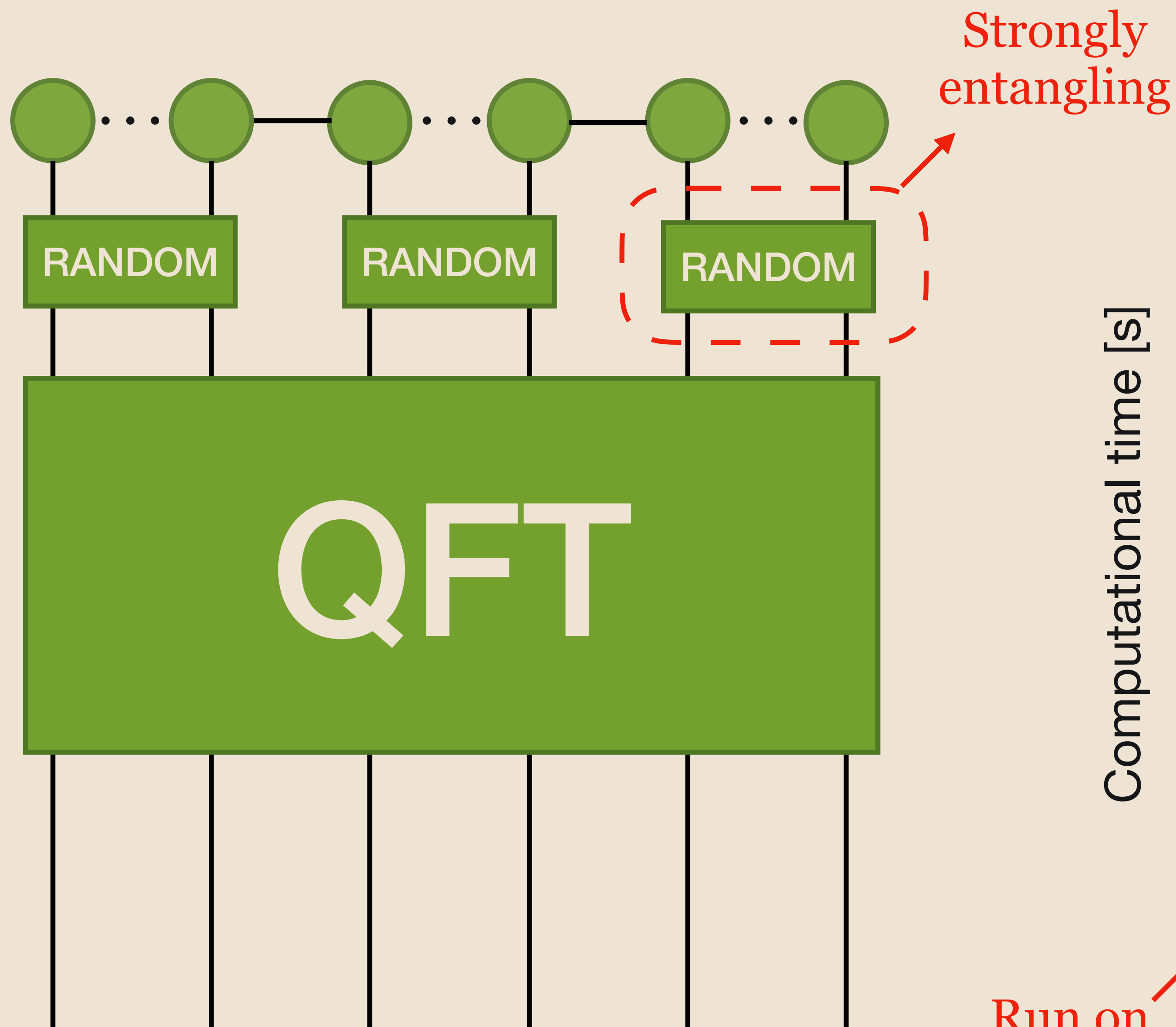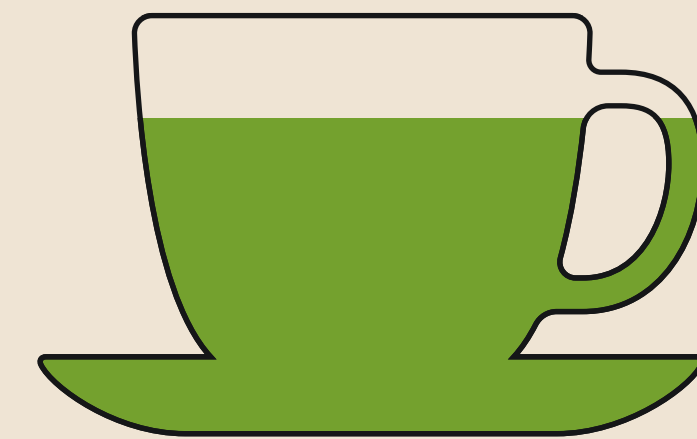# **Benchmarks**

# Benchmarks

# Benchmarks



Strongly entangling

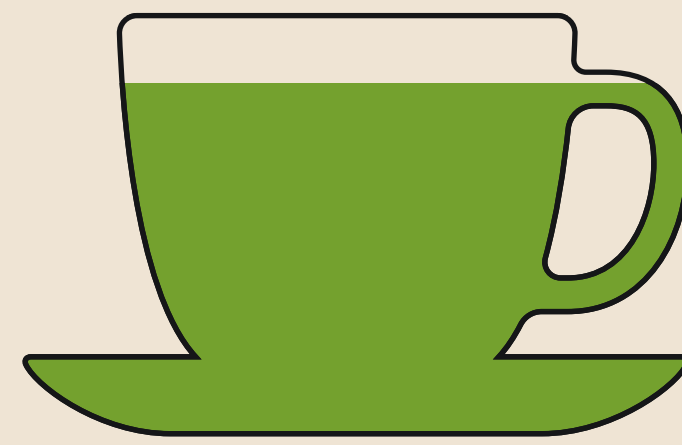RANDOM RANDOM RANDOM

QFT

$\chi_{max}$ reached

$n = 100, \chi_{max} = 1024, 16$ threads

Computational time [s]

Entangling block size

- qiskit
- fortran
- numpy
- Wall time

Run on Galileo100

# Applications

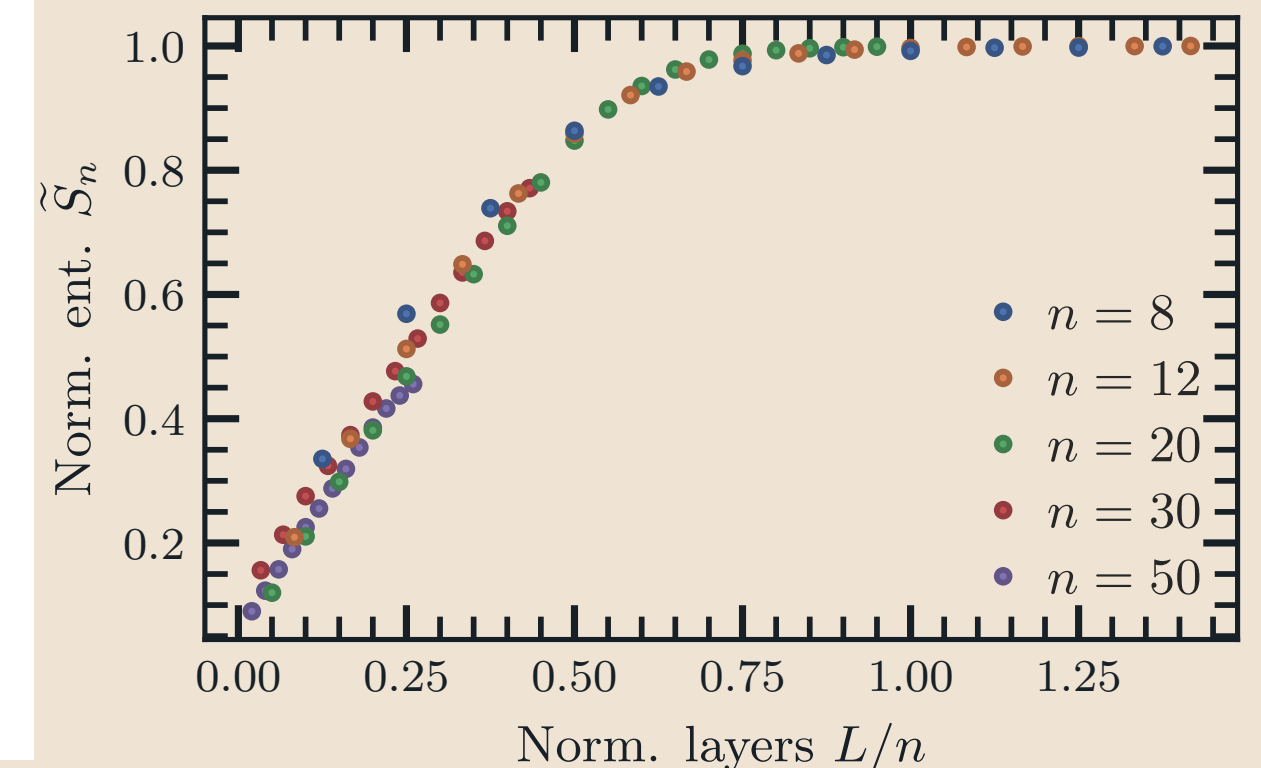**Entanglement entropy production in QNN**
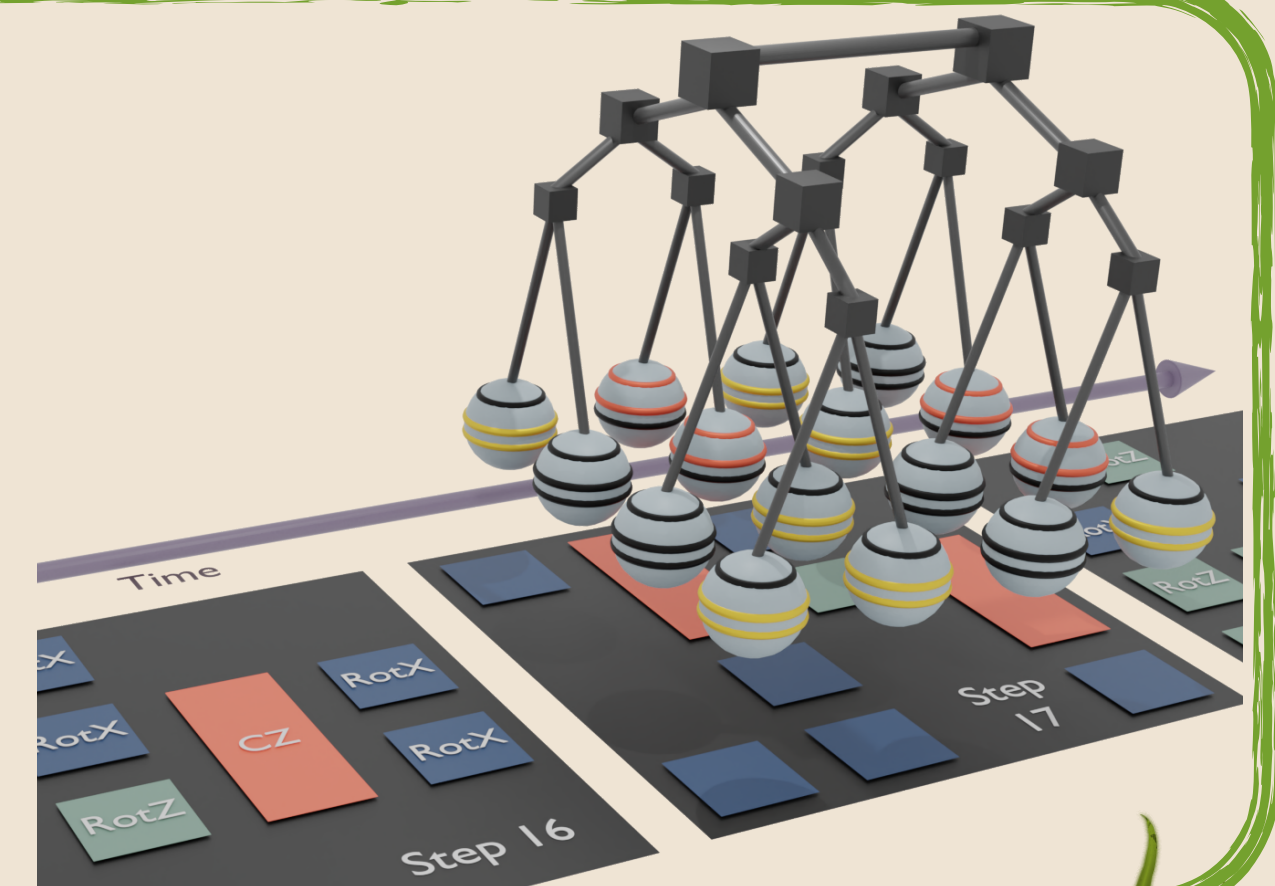Ballarin, Marco, et al. arXiv:2206.02474
- Simulations up to 50 qubits
- Bond dimension of 4096
- 11h of runtime on Galileo100
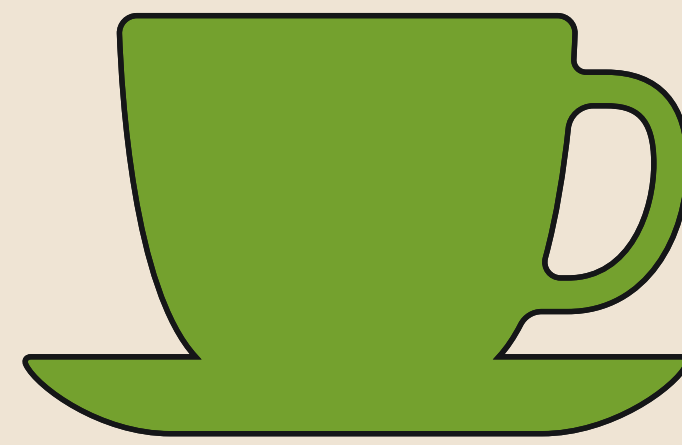


**Ab initio two-dimensional digital twin for quantum computer**
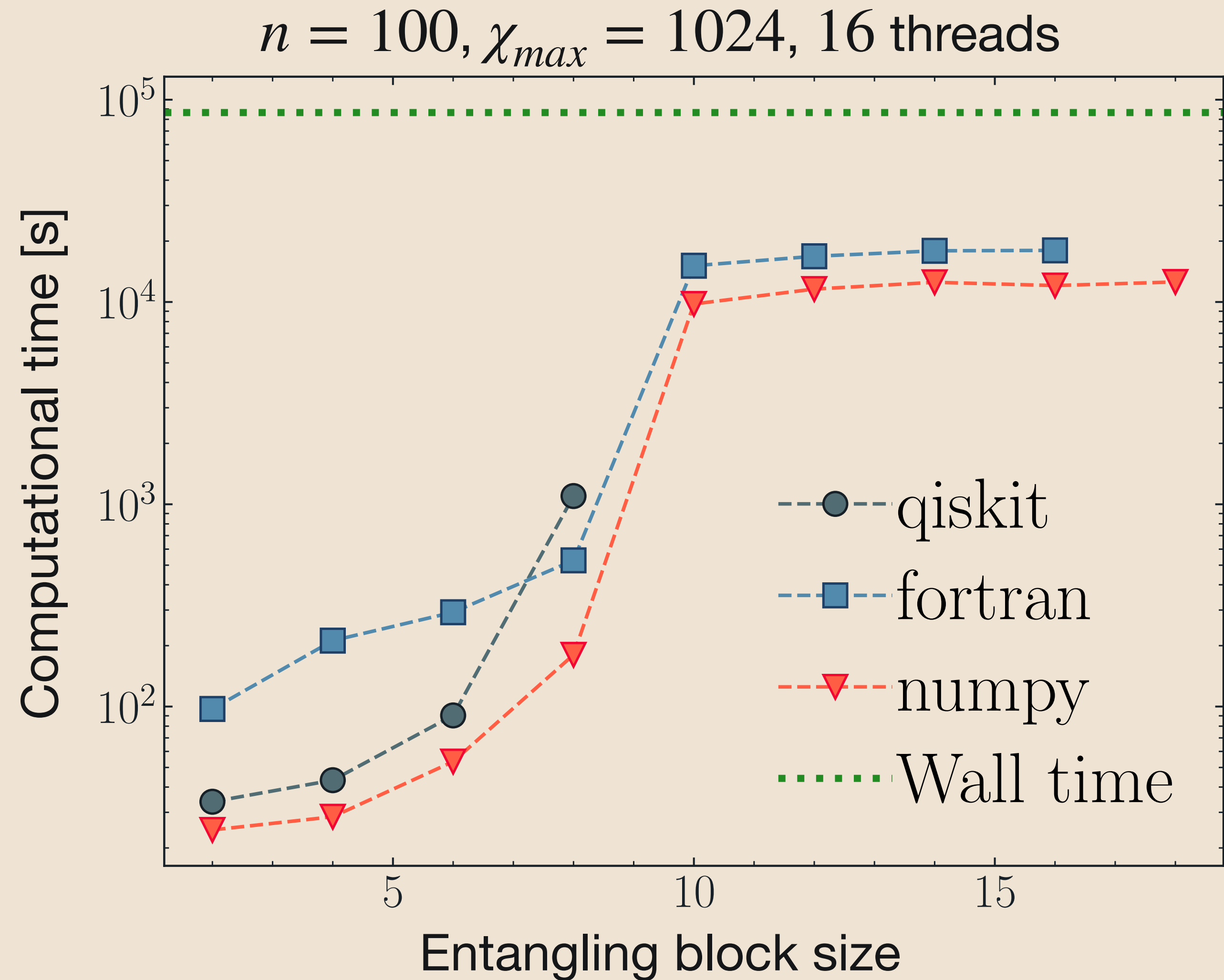Jaschke, Daniel, et al. arXiv:2210.03763
- Use of the unbiased sampling
- Quantum matcha tea simulations used as target state to compute the fidelity of a simulation with crosstalk
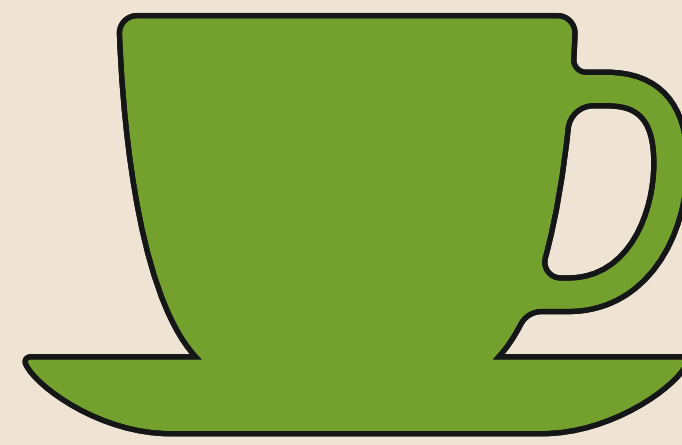
# Conclusions

MPS simulations are not limited by the number of qubits but by the entanglement



$n = 100, \chi_{max} = 1024, 16$ threads

Computational time [s]

- qiskit
- fortran
- numpy
- Wall time

Entangling block size

# Conclusions

MPS simulations are not limited by the number of qubits but by the entanglement

Easy-to-use python frontend and fast HPC-ready backend (Both GPU and CPU)



$n = 100, \chi_{max} = 1024, 16$ threads

qiskit
fortran
numpy
Wall time

Computational time [s]
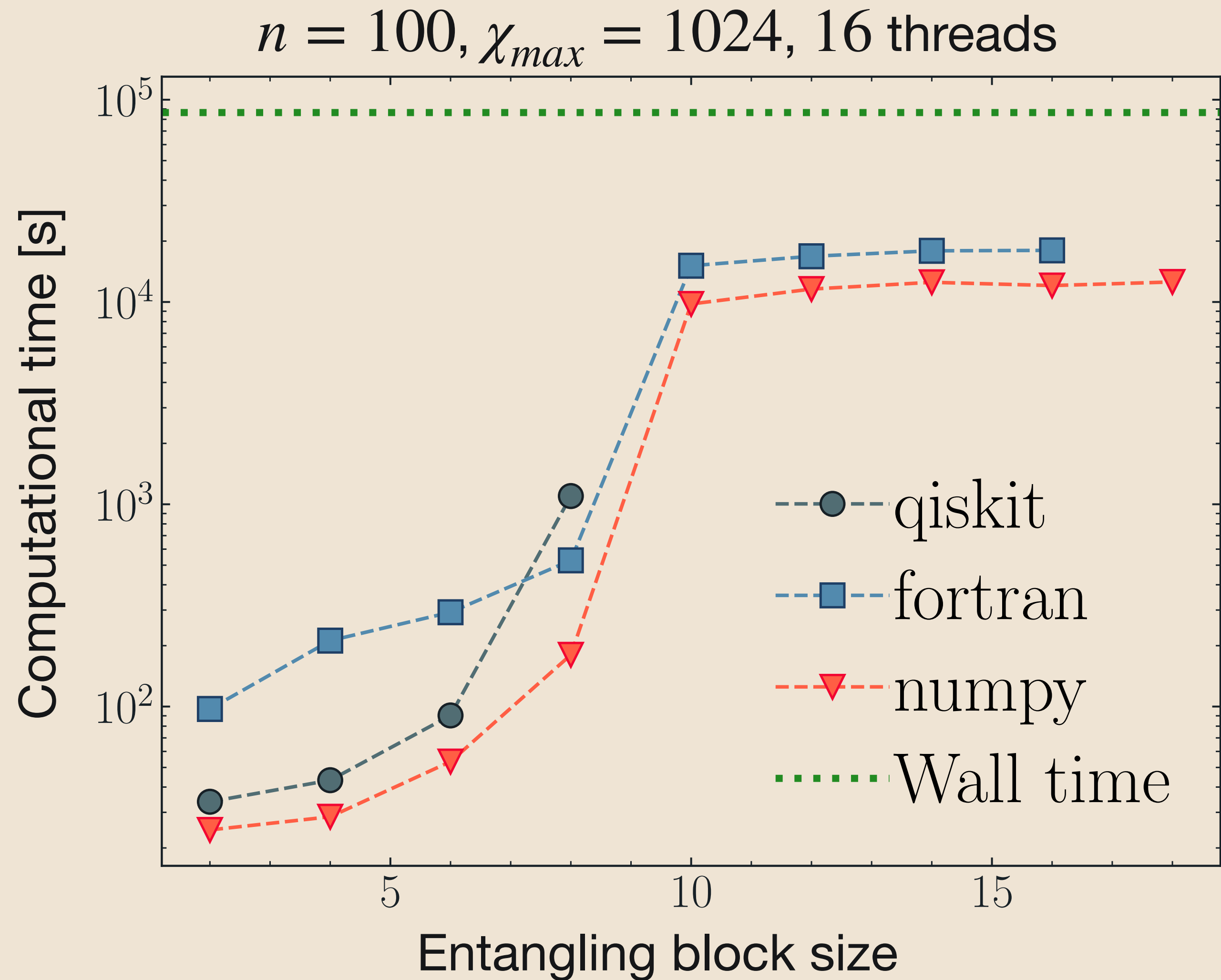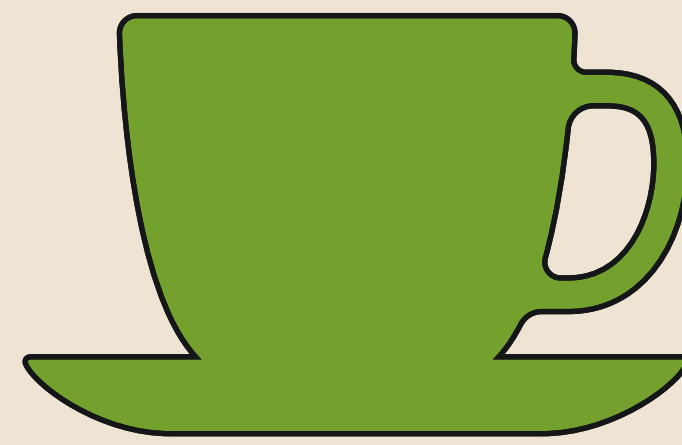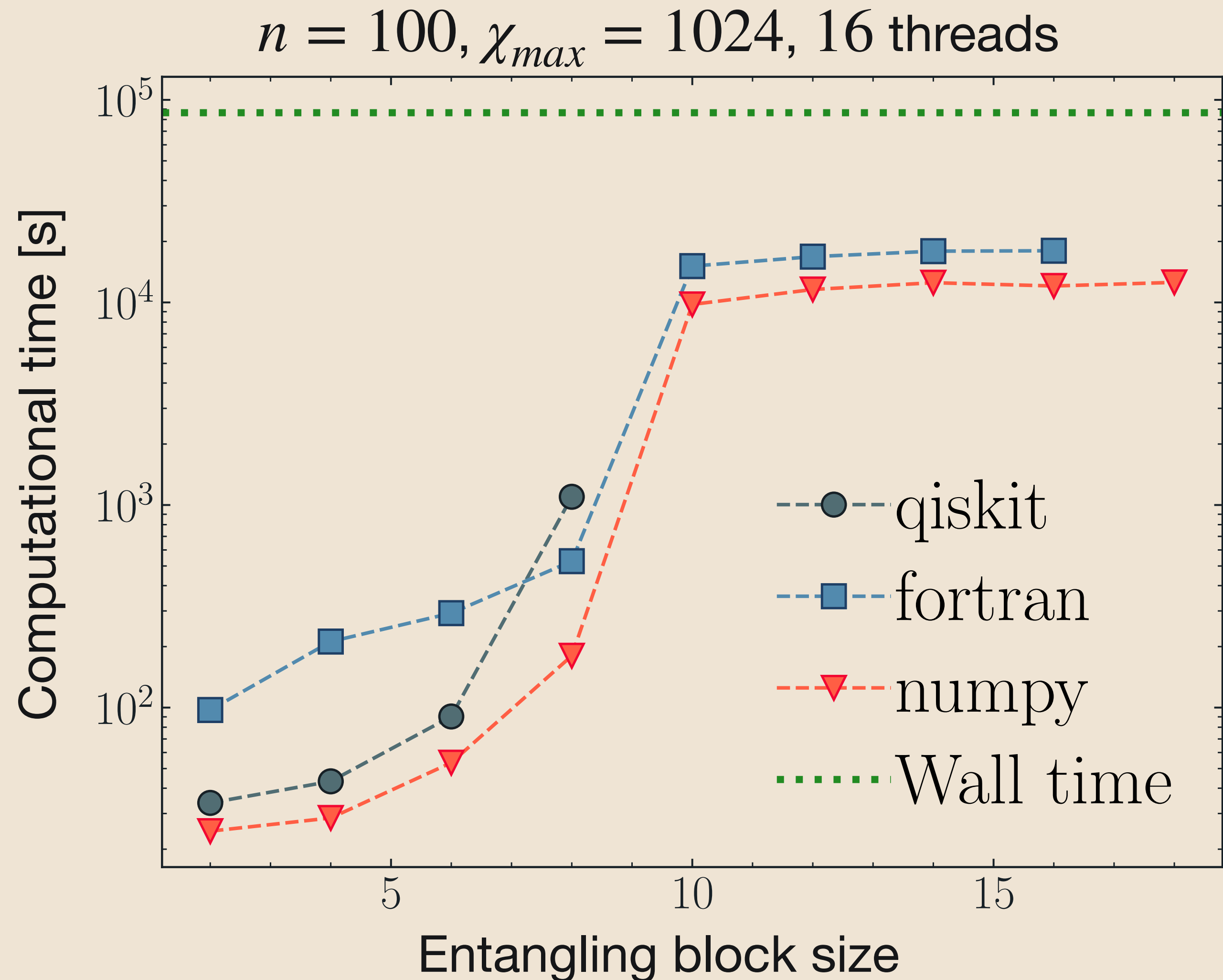
Entangling block size

# Conclusions

MPS simulations are not limited by the number of qubits but by the entanglement

Easy-to-use python frontend and fast HPC-ready backend (Both GPU and CPU)

Error analysis tools and efficient computations of observables optimised for the MPS representation

$n = 100, \chi_{max} = 1024, 16$ threads



Computational time [s]

Entangling block size

- qiskit
- fortran
- numpy
- Wall time
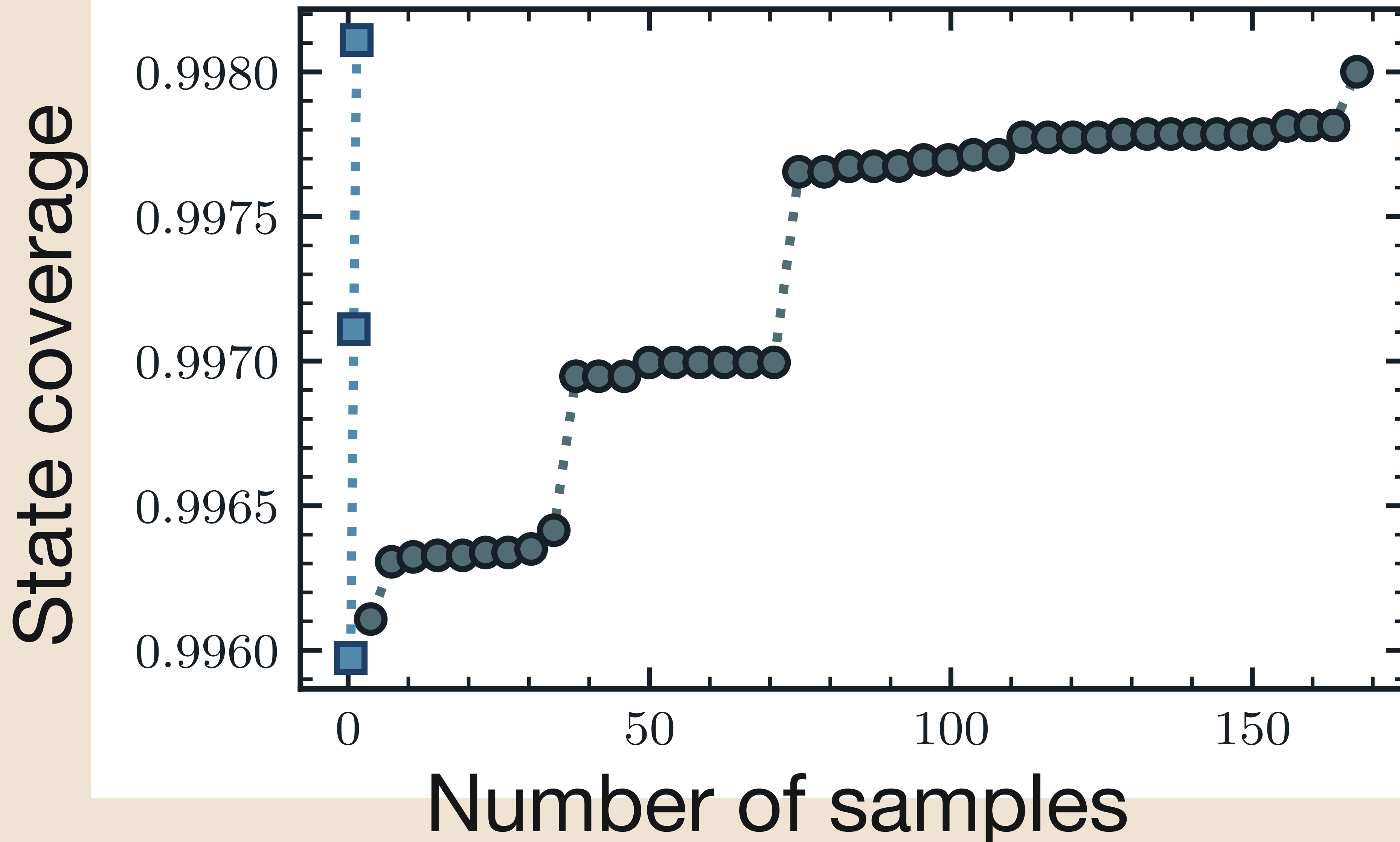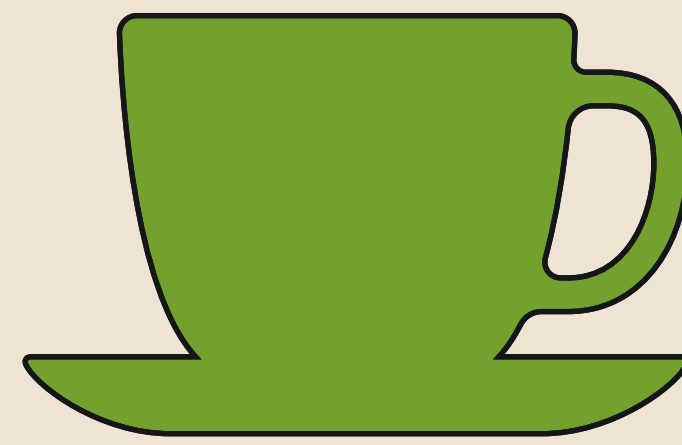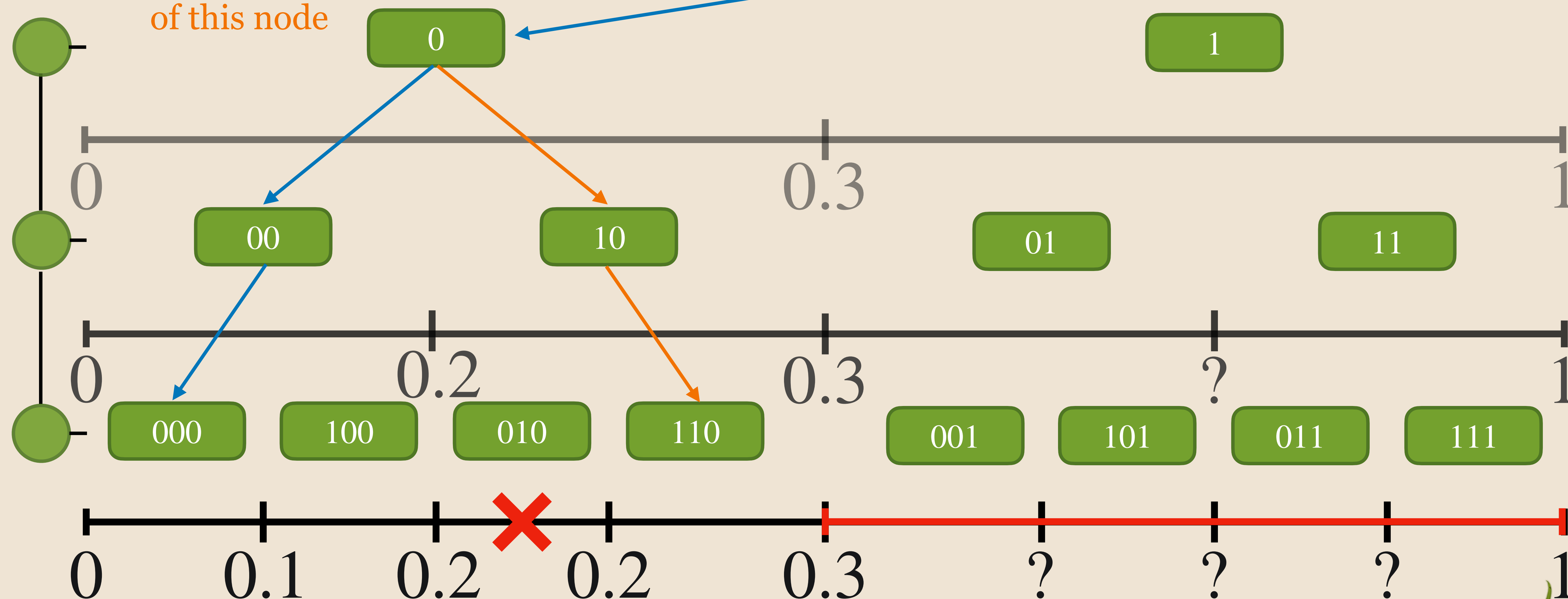
# Efficient sampling of final state

# Efficient sampling of final state

Reuse the computation of this node

Sample random number $n = 0.05, 0.29$



We know which states we did not sample and can sample only here in second round

15